

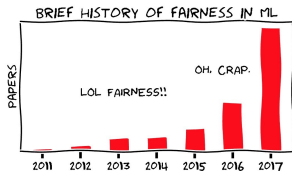
STA414/2104

Statistical Methods for Machine Learning II

Murat A. Erdogdu

Department of Computer Science
Department of Statistical Sciences

Lecture 11



Announcements

- This is the last lecture!
 - ▶ Algorithmic fairness.
 - ▶ Final exam review.
- Please fill out the course evaluation form. Final response rate will be shared with class (currently 68%).

Overview

- As ML starts to be applied to critical applications involving humans, the field is wrestling with the societal impacts
 - ▶ **Security:** what if an attacker tries to poison the training data, fool the system with malicious inputs, “steal” the model, etc.?

Overview

- As ML starts to be applied to critical applications involving humans, the field is wrestling with the societal impacts
 - ▶ **Security:** what if an attacker tries to poison the training data, fool the system with malicious inputs, “steal” the model, etc.?
 - ▶ **Privacy:** avoid leaking (much) information about the data the system was trained on (e.g. medical diagnosis)

Overview

- As ML starts to be applied to critical applications involving humans, the field is wrestling with the societal impacts
 - ▶ **Security:** what if an attacker tries to poison the training data, fool the system with malicious inputs, “steal” the model, etc.?
 - ▶ **Privacy:** avoid leaking (much) information about the data the system was trained on (e.g. medical diagnosis)
 - ▶ **Fairness:** ensure that the system doesn’t somehow disadvantage particular individuals or groups

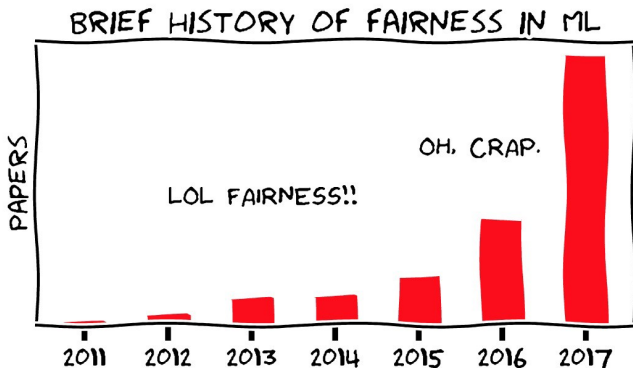
Overview

- As ML starts to be applied to critical applications involving humans, the field is wrestling with the societal impacts
 - ▶ **Security:** what if an attacker tries to poison the training data, fool the system with malicious inputs, “steal” the model, etc.?
 - ▶ **Privacy:** avoid leaking (much) information about the data the system was trained on (e.g. medical diagnosis)
 - ▶ **Fairness:** ensure that the system doesn’t somehow disadvantage particular individuals or groups
 - ▶ **Transparency:** be able to understand why one decision was made rather than another

Overview

- As ML starts to be applied to critical applications involving humans, the field is wrestling with the societal impacts
 - ▶ **Security:** what if an attacker tries to poison the training data, fool the system with malicious inputs, “steal” the model, etc.?
 - ▶ **Privacy:** avoid leaking (much) information about the data the system was trained on (e.g. medical diagnosis)
 - ▶ **Fairness:** ensure that the system doesn’t somehow disadvantage particular individuals or groups
 - ▶ **Transparency:** be able to understand why one decision was made rather than another
 - ▶ **Accountability:** an outside auditor should be able to verify that the system is functioning as intended
- If some of these definitions sound vague, that’s because formalizing them is the main challenge!

Overview: Fairness



Credit: Moritz Hardt

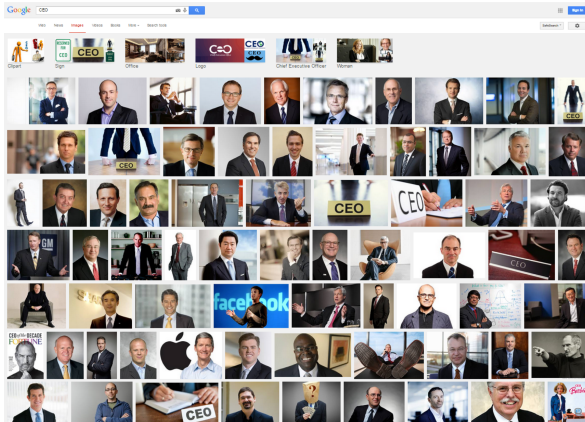
Overview: Fairness

FAIRNESS IN AUTOMATED DECISIONS



Overview: Fairness

SUBTLER BIAS



Overview: Fairness

The image shows two screenshots of the Google Translate interface. The top screenshot shows the translation of English text to Turkish. The bottom screenshot shows the translation of Turkish text to English.

Top Screenshot:

- Language pair: English to Turkish.
- Input text: "She is a doctor. He is a nurse."
- Output text: "O bir doktor. O bir hemşire."
- Character count: 31/5000.

Bottom Screenshot:

- Language pair: Turkish to English.
- Input text: "O bir doktor. O bir hemşire"
- Output text: "He is a doctor. She is a nurse ✓"
- Character count: 28/5000.

Turkish has gender neutral pronouns

Overview: Fairness

- This lecture: algorithmic fairness
- Goal: identify and mitigate **bias** in ML-based decision making, in all aspects of the pipeline
- Sources of bias/discrimination
 - ▶ Data
 - ▶ Imbalanced/impoverished data
 - ▶ Labeled data imbalance
 - ▶ Labeled data incorrect / noisy
 - ▶ Model
 - ▶ ML prediction error imbalanced
 - ▶ Compound injustices

Overview: Fairness

- This lecture: algorithmic fairness
- Goal: identify and mitigate **bias** in ML-based decision making, in all aspects of the pipeline
- Sources of bias/discrimination
 - ▶ Data
 - ▶ Imbalanced/impoverished data
 - ▶ Labeled data imbalance
 - ▶ Labeled data incorrect / noisy
 - ▶ Model
 - ▶ ML prediction error imbalanced
 - ▶ Compound injustices
- Important: Algorithmic fairness does not imply real fairness!

An Example from FairML Book

- PredPol is a predictive policing algorithm.
- By applying PredPol to data derived from Oakland police records, researchers found that Black people would be targeted for predictive policing of drug crimes at roughly twice the rate of whites, even though the two groups have roughly equal rates of drug use.
- Their analysis showed that this initial bias would be amplified by a feedback loop, with policing increasingly concentrated on targeted areas.
- This is despite the fact that the PredPol algorithm does not explicitly take demographics into account.

Learning Fair Representations

- In a classification problem, there may be sensitive attributes which cause bias.
- A naïve attempt: simply don't use the sensitive feature.

Learning Fair Representations

- In a classification problem, there may be sensitive attributes which cause bias.
- A naïve attempt: simply don't use the sensitive feature.
 - ▶ Problem: the algorithm implicitly learns to predict the sensitive feature from other features (e.g. race from zip code)
- Another idea: limit the algorithm to a small set of features you're pretty sure are safe and task-relevant
 - ▶ This is the conservative approach, and commonly used for both human and machine decision making

Learning Fair Representations

- In a classification problem, there may be sensitive attributes which cause bias.
- A naïve attempt: simply don't use the sensitive feature.
 - ▶ Problem: the algorithm implicitly learns to predict the sensitive feature from other features (e.g. race from zip code)
- Another idea: limit the algorithm to a small set of features you're pretty sure are safe and task-relevant
 - ▶ This is the conservative approach, and commonly used for both human and machine decision making
 - ▶ But removing features hurts the classification accuracy. Maybe we can make more accurate decisions if we include more features and somehow enforce fairness algorithmically?
- Can we learn fair representations, which can make accurate classifications without implicitly using the sensitive attribute?

Overview: Fairness

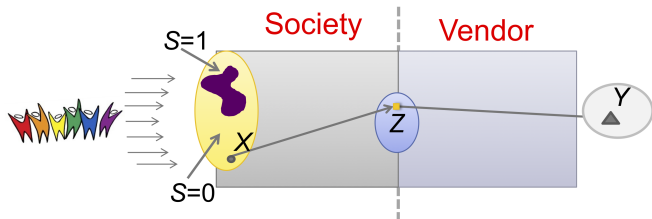
- Notation
 - ▶ $X \in \mathbb{R}^D$: input to classifier
 - ▶ $S \in \{0, 1\}$: belongs to protected group (age, gender, race, etc.)
 - ▶ $Z \in \{1, 2, \dots, K\}$: latent representation (code)
 - ▶ $T \in \{0, 1\}$: true label
 - ▶ $Y \in [0, 1]$: prediction ($p(T = 1 | X)$)
- We use capital letters to emphasize that these are random variables.

Fairness Criteria

- $X \perp\!\!\!\perp Y$ means X and Y are independent
- Most common way to define fair classification is to require some invariance with respect to the sensitive attribute
 - ▶ Demographic parity: $Y \perp\!\!\!\perp S$
 - ▶ Equalized odds: $Y \perp\!\!\!\perp S \mid T$
 - ▶ Equal opportunity: $Y \perp\!\!\!\perp S \mid T = t$, for a fixed t
 - ▶ Equal (weak) calibration: $T \perp\!\!\!\perp S \mid Y$
 - ▶ Equal (strong) calibration: $T \perp\!\!\!\perp S \mid Y$ and $Y = p(T = 1)$
 - ▶ Fair subgroup accuracy: $\mathbb{1}[T = Y] \perp\!\!\!\perp S$
- Many of these definitions are incompatible!

Learning Fair Representations

- Idea: separate the responsibilities of the (trusted) society and (untrusted) vendor



- Goal: find a representation Z that removes any information about the sensitive attribute
- Then the vendor can do whatever they want!

Learning Fair Representations

Desiderata for the representation:

- Retain information about $X \Rightarrow$ high mutual information between X and Z
- Obfuscate $S \Rightarrow$ low mutual information between S and Z
- Allow high classification accuracy \Rightarrow high mutual information between T and Z

Learning Fair Representations

First approach: Zemel et al., 2013, "Learning fair representations"

- We want to train a fair classifier.
- We can solve the following problem

$$\mathcal{L}_{\text{total}} = \lambda_r \mathcal{L}_{\text{reconst}} + \lambda_p \mathcal{L}_{\text{pred}} + \lambda_d \mathcal{L}_{\text{discrim}}$$

where λ_r , λ_p , and λ_d are hyperparameters governing the trade-off between losses.

- ▶ Each loss has a different job to be defined next.

Learning Fair Representations

- Let Z be a discrete code or representation (like K-means)

Learning Fair Representations

- Let Z be a discrete code or representation (like K-means)
- Determine Z based on distance to (the cluster center in K-means)

$$r_{ik} = p(Z = k | x_i) \propto \exp(-\beta \|x_i - \mu_k\|^2),$$

where $\beta > 0$ is a constant, and μ_k is a prototype for the cluster.

- Need to fit the prototypes μ_k . They are unknown.

Learning Fair Representations

- Let Z be a discrete code or representation (like K-means)
- Determine Z based on distance to (the cluster center in K-means)

$$r_{ik} = p(Z = k | x_i) \propto \exp(-\beta \|x_i - \mu_k\|^2),$$

where $\beta > 0$ is a constant, and μ_k is a prototype for the cluster.

- Need to fit the prototypes μ_k . They are unknown.
- Similar to EM update, we let the reconstruction be

$$\tilde{x}_i = \sum_{k=1}^K r_{ik} \mu_k$$

and enforce that $x_i \approx \tilde{x}_i$ by minimizing

$$\mathcal{L}_{\text{reconst}} = \frac{1}{N} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2.$$

Learning Fair Representations

- Remember, we want to train a fair **classifier**.
- We predict using a linear function of $\mathbf{r}_i = [r_{i1}, r_{i2}, \dots, r_{ik}]^\top$.

$$y_i = \sigma(\mathbf{w}^\top \mathbf{r}_i) = p(t_i | \mathbf{x}_i)$$

- Need to fit weights \mathbf{w} . They are unknown.
- Loss:

Learning Fair Representations

- Remember, we want to train a fair **classifier**.
- We predict using a linear function of $\mathbf{r}_i = [r_{i1}, r_{i2}, \dots, r_{ik}]^\top$.

$$y_i = \sigma(\mathbf{w}^\top \mathbf{r}_i) = p(t_i | \mathbf{x}_i)$$

- Need to fit weights \mathbf{w} . They are unknown.
- Loss: we can use cross-entropy

$$L_{\text{CE}}(y_i, t_i) = -t_i \log y_i - (1 - t_i) \log(1 - y_i)$$

Learning Fair Representations

- Next, enforce a fairness constraint:

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i) \right|.$$

- $N_0 = \#\{i : s_i = 0\}$, $N_1 = \#\{i : s_i = 1\}$ and $N_0 + N_1 = N$.
- We show this enforces **demographic parity**.

Learning Fair Representations

- Enforce **demographic parity** by obfuscating S :

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i) \right|,$$

- $N_0 = \#\{i : s_i = 0\}$, $N_1 = \#\{i : s_i = 1\}$ and $N_0 + N_1 = N$.

Learning Fair Representations

- Enforce **demographic parity** by obfuscating S :

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i) \right|,$$

- $N_0 = \#\{i : s_i = 0\}$, $N_1 = \#\{i : s_i = 1\}$ and $N_0 + N_1 = N$.
- If the above discrimination loss is $\mathcal{L}_{\text{discrim}} = 0$, we have **LHS=RHS** for all $k = 1, 2, \dots, K$. Therefore,

Learning Fair Representations

- Enforce **demographic parity** by obfuscating S :

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i) \right|,$$

- $N_0 = \#\{i : s_i = 0\}$, $N_1 = \#\{i : s_i = 1\}$ and $N_0 + N_1 = N$.
- If the above discrimination loss is $\mathcal{L}_{\text{discrim}} = 0$, we have **LHS=RHS** for all $k = 1, 2, \dots, K$. Therefore,

$$p(Y = 1 | S = 1) = \sum_k p(Y = 1 | Z = k) p(Z = k | S = 1)$$

Learning Fair Representations

- Enforce **demographic parity** by obfuscating S :

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i) \right|,$$

- $N_0 = \#\{i : s_i = 0\}$, $N_1 = \#\{i : s_i = 1\}$ and $N_0 + N_1 = N$.
- If the above discrimination loss is $\mathcal{L}_{\text{discrim}} = 0$, we have **LHS=RHS** for all $k = 1, 2, \dots, K$. Therefore,

$$\begin{aligned} p(Y = 1 | S = 1) &= \sum_k p(Y = 1 | Z = k) p(Z = k | S = 1) \\ &\approx \sum_k p(Y = 1 | Z = k) \mathbb{E}_x \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i = 1) \end{aligned}$$

Learning Fair Representations

- Enforce **demographic parity** by obfuscating S :

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i) \right|,$$

- $N_0 = \#\{i : s_i = 0\}$, $N_1 = \#\{i : s_i = 1\}$ and $N_0 + N_1 = N$.
- If the above discrimination loss is $\mathcal{L}_{\text{discrim}} = 0$, we have **LHS=RHS** for all $k = 1, 2, \dots, K$. Therefore,

$$\begin{aligned} p(Y = 1 | S = 1) &= \sum_k p(Y = 1 | Z = k) p(Z = k | S = 1) \\ &\approx \sum_k p(Y = 1 | Z = k) \mathbb{E}_x \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i = 1) \\ &= \sum_k p(Y = 1 | Z = k) \mathbb{E}_x \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i = 0) \end{aligned}$$

Learning Fair Representations

- Enforce **demographic parity** by obfuscating S :

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i) \right|,$$

- $N_0 = \#\{i : s_i = 0\}$, $N_1 = \#\{i : s_i = 1\}$ and $N_0 + N_1 = N$.
- If the above discrimination loss is $\mathcal{L}_{\text{discrim}} = 0$, we have **LHS=RHS** for all $k = 1, 2, \dots, K$. Therefore,

$$\begin{aligned} p(Y = 1 | S = 1) &= \sum_k p(Y = 1 | Z = k) p(Z = k | S = 1) \\ &\approx \sum_k p(Y = 1 | Z = k) \mathbb{E}_x \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i = 1) \\ &= \sum_k p(Y = 1 | Z = k) \mathbb{E}_x \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i = 0) \\ &\approx \sum_k p(Y = 1 | Z = k) p(Z = k | S = 0) \end{aligned}$$

Learning Fair Representations

- Enforce **demographic parity** by obfuscating S :

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i) \right|,$$

- $N_0 = \#\{i : s_i = 0\}$, $N_1 = \#\{i : s_i = 1\}$ and $N_0 + N_1 = N$.
- If the above discrimination loss is $\mathcal{L}_{\text{discrim}} = 0$, we have **LHS=RHS** for all $k = 1, 2, \dots, K$. Therefore,

$$\begin{aligned} p(Y = 1 | S = 1) &= \sum_k p(Y = 1 | Z = k) p(Z = k | S = 1) \\ &\approx \sum_k p(Y = 1 | Z = k) \mathbb{E}_x \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i, s_i = 1) \\ &= \sum_k p(Y = 1 | Z = k) \mathbb{E}_x \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i, s_i = 0) \\ &\approx \sum_k p(Y = 1 | Z = k) p(Z = k | S = 0) \\ &= p(Y = 1 | S = 0) \quad \text{demographic parity} \end{aligned}$$

Learning Fair Representations

- We want to retain information about X : $x_i \approx \tilde{x}_i$; penalize reconstruction error

$$\mathcal{L}_{\text{reconst}} = \frac{1}{N} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2$$

- Predict accurately: cross-entropy loss

$$\mathcal{L}_{\text{pred}} = \frac{1}{N} \sum_{i=1}^N -t_i \log y_i - (1 - t_i) \log(1 - y_i)$$

- Obfuscate S :

$$\mathcal{L}_{\text{discrim}} = \frac{1}{K} \sum_{k=1}^K \left| \frac{1}{N_0} \sum_{i:s_i=0} p(Z = k | x_i) - \frac{1}{N_1} \sum_{i:s_i=1} p(Z = k | x_i) \right|.$$

Learning Fair Representations

- We can solve the following problem

$$\mathcal{L}_{\text{total}}(\{\mu_k\}_{k=1}^K, \mathbf{w}) = \lambda_r \mathcal{L}_{\text{reconst}} + \lambda_p \mathcal{L}_{\text{pred}} + \lambda_d \mathcal{L}_{\text{discrim}}$$

where λ_r , λ_p , and λ_d are hyperparameters governing the trade-off between losses.

- We can find the optimal parameter $\{\mu_k\}_{k=1}^K, \mathbf{w}$ using an alternating optimization method.

Learning Fair Representations

Datasets

1. German Credit

Task: classify individual as good or bad credit risk

Sensitive feature: Age

2. Adult Income

Size: 45,222 instances, 14 attributes

Task: predict whether or not annual income > 50K

Sensitive feature: Gender

3. Heritage Health

Size: 147,473 instances, 139 attributes

Task: predict whether patient spends any nights in hospital

Sensitive feature: Age

Learning Fair Representations

Metrics

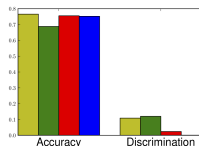
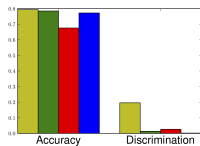
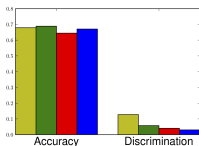
- Classification accuracy
- Discrimination: measuring the difference in proportion of positive classification of individuals in the protected or unprotected groups.

$$\left| \frac{\sum_{i:s_i=1}^N y_i}{N_1} - \frac{\sum_{i:s_i=0}^N y_i}{N_0} \right|$$

German

Adult

Health



Blue = theirs, others: logistic reg (LR), naive Bayes, regularized LR

Individual Fairness

- The work on fair representations was geared towards group fairness
- Another notion of fairness is individual level: ensuring that similar individuals are treated similarly by the algorithm
 - ▶ This depends heavily on the notion of “similar”.
- One way to define similarity is in terms of the “true label” T (e.g. whether this individual is in fact likely to repay their loan)
 - ▶ Can you think of a problem with this definition?

Individual Fairness

- The work on fair representations was geared towards group fairness
- Another notion of fairness is individual level: ensuring that similar individuals are treated similarly by the algorithm
 - ▶ This depends heavily on the notion of “similar”.
- One way to define similarity is in terms of the “true label” T (e.g. whether this individual is in fact likely to repay their loan)
 - ▶ Can you think of a problem with this definition?
 - ▶ The label may itself be biased
 - ▶ if based on human judgments
 - ▶ if, e.g., societal biases make it harder for one group to pay off their loans
 - ▶ Keep in mind that you'd need to carefully consider the assumptions when applying one of these methods!

Equal Opportunity

- Example: predict whether an individual is likely to repay their loan
- Two notions of individual fairness:
 - ▶ **Equalized odds**: equal true positive and false positive rates

$$p(Y = 1 | S = 0, T = t) = p(Y = 1 | S = 1, T = t) \quad \text{for } t \in \{0, 1\}$$

Equal Opportunity

- Example: predict whether an individual is likely to repay their loan
- Two notions of individual fairness:
 - ▶ **Equalized odds:** equal true positive and false positive rates

$$p(Y = 1 | S = 0, T = t) = p(Y = 1 | S = 1, T = t) \quad \text{for } t \in \{0, 1\}$$

- ▶ **Equal opportunity:** equal true positive rates

$$p(Y = 1 | S = 0, T = 1) = p(Y = 1 | S = 1, T = 1)$$

Fairness Summary

- Fairness is a challenging issue to address
 - ▶ Not something you can just measure on a validation set
 - ▶ Philosophers and lawyers have been trying to define it for thousands of years
 - ▶ Different notions are incompatible. Need to carefully consider the particular problem.
 - ▶ individual vs. group
- Explosion of interest in ML over the last few years
- Conference on Fairness, Accountability, and Transparency (FAT*)
- New textbook: <https://fairmlbook.org/>

Final Exam

- Final exam is on 4/20, at 9am EST.
- Covers all lectures except the one on Algorithmic Fairness.
- It doesn't cover material only covered in suggested reading.
- Similar difficulty to midterm.
- Practice final exam will be posted.

Basic ML Terminology

The final exam will be on the entire course; however, it will be more weighted towards post-midterm material. For pre-midterm material, refer to the midterm review slides on the website.

- Generalization, overfitting, underfitting
- Regression/Classification
- Bias-Variance decomp
- Training, test, validation
- Neural Network, backprop
- Cross-entropy, softmax
- MLE, MAP, prior, posterior
- PCA, Dimension Reduction, projection on subspace
- Matrix factorization, EVD
- Generative vs. Discriminative Methods, Bayes rule
- Trees, Gaussian Discriminant Analysis
- SVMs, Ensembles: Bagging
- k-Means, EM algorithm
- Value function, Q-learning, Bellman equations

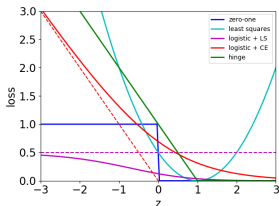
Lec 7: Support Vector Machines

Previously, we saw loss functions as relaxations to 0-1 loss. For SVMs, we consider a different relaxation: **Hinge loss**

$$\text{0-1 loss: } \mathcal{L}_{0-1}(z, t) = \mathbb{I}\{\text{sign}(z) \neq t\}$$

$$\text{Hinge loss: } \mathcal{L}_H(z, t) = \max\{0, 1 - zt\}$$

Classification w/ Different losses



Lec 7: Support Vector Machines

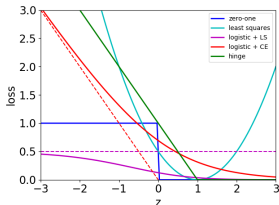
Previously, we saw loss functions as relaxations to 0-1 loss. For SVMs, we consider a different relaxation: **Hinge loss**

$$\text{0-1 loss: } \mathcal{L}_{0-1}(z, t) = \mathbb{I}\{\text{sign}(z) \neq t\}$$

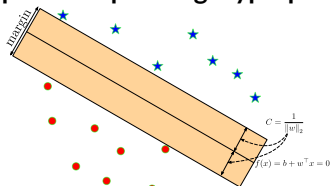
$$\text{Hinge loss: } \mathcal{L}_H(z, t) = \max\{0, 1 - zt\}$$

$$\text{SVM loss: } \min_{w, b} \sum_{i=1}^N \max\{0, 1 - t_i y_i(w, b)\} + \frac{\lambda}{2} \|w\|_2^2$$

Classification w/ Different losses



What does SVM do?: Optimal Separating Hyperplane



Lec 7: SVMs and Maximizing Margin

$$\begin{aligned} & \max_{w,b} C \\ \text{s.t. } & \frac{t_i(w^\top x_i + b)}{\|w\|_2} \geq C \quad i = 1, \dots, N \end{aligned}$$

Lec 7: SVMs and Maximizing Margin

$$\begin{aligned} & \max_{w,b} C \\ \text{s.t.} & \frac{t_i(w^\top x_i + b)}{\|w\|_2} \geq C \quad i = 1, \dots, N \\ & \|w\|_2 = \frac{1}{C} \end{aligned}$$

Lec 7: SVMs and Maximizing Margin

$$\begin{aligned} & \max_{w,b} C \\ & \text{s.t. } \frac{t_i(w^\top x_i + b)}{\|w\|_2} \geq C \quad i = 1, \dots, N \\ & \quad \|w\|_2 = \frac{1}{C} \end{aligned}$$

Note that if $\|w\|_2 = \frac{1}{C}$ then:

$$\underbrace{\frac{t_i(w^\top x_i + b)}{\|w\|_2} \geq \frac{1}{\|w\|_2}}_{\text{geometric margin constraint}} \iff \underbrace{t_i(w^\top x_i + b) \geq 1}_{\text{algebraic margin constraint}}$$

Lec 7: SVMs and Maximizing Margin

$$\begin{aligned} & \max_{w,b} C \\ & \text{s.t. } \frac{t_i(w^\top x_i + b)}{\|w\|_2} \geq C \quad i = 1, \dots, N \\ & \quad \|w\|_2 = \frac{1}{C} \end{aligned}$$

Note that if $\|w\|_2 = \frac{1}{C}$ then:

$$\underbrace{\frac{t_i(w^\top x_i + b)}{\|w\|_2} \geq \frac{1}{\|w\|_2}}_{\text{geometric margin constraint}} \iff \underbrace{t_i(w^\top x_i + b) \geq 1}_{\text{algebraic margin constraint}}$$

Plugging in $C = \frac{1}{\|w\|_2}$, equivalent optimization objective:

$$\begin{aligned} & \min \|w\|_2^2 \\ & \text{s.t. } t_i(w^\top x_i + b) \geq 1 \quad i = 1, \dots, N \end{aligned}$$

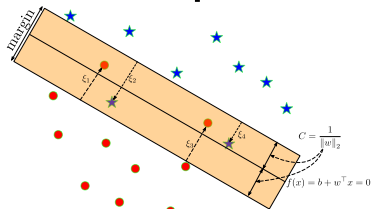
Lec 7: SVMs w/ Non-Separable Data

- Soft margin constraint:

$$\frac{t_i(w^\top x_i + b)}{\|w\|_2} \geq C(1 - \xi_i),$$

for $\xi_i \geq 0$.

- Penalize $\sum_i \xi_i$



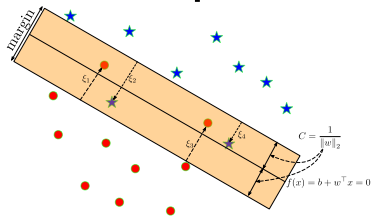
Lec 7: SVMs w/ Non-Separable Data

- **Soft margin constraint:**

$$\frac{t_i(w^\top x_i + b)}{\|w\|_2} \geq C(1 - \xi_i),$$

for $\xi_i \geq 0$.

- Penalize $\sum_i \xi_i$



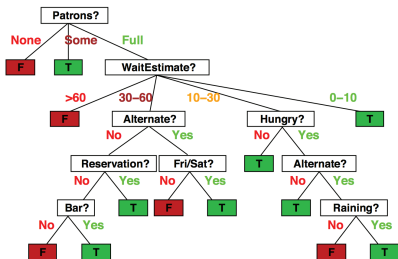
The **Soft-margin SVM** objective becomes:

$$\begin{aligned} \min_{w, b, \xi} \quad & \frac{1}{2} \|w\|_2^2 + \gamma \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & t_i(w^\top x_i + b) \geq 1 - \xi_i \quad i = 1, \dots, N \\ & \xi_i \geq 0 \quad i = 1, \dots, N \end{aligned}$$

If we write $y_i(w, b) = w^\top x_i + b$, then it can be written as

$$\min_{w, b} \sum_{i=1}^N \max\{0, 1 - t_i y_i(w, b)\} + \frac{1}{2\gamma} \|w\|_2^2$$

Lec 7: Decision Trees

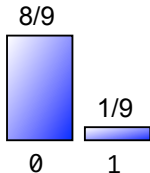


- **Internal nodes** test **attributes**
- **Branching** is determined by **attribute value**
- **Leaf nodes** are **outputs** (predictions)

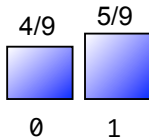
Lec 7: Quantifying Uncertainty

Entropy is a measure of expected “surprise”: How uncertain are we of the value of a draw from this distribution?

$$H(X) = -\mathbb{E}_{X \sim p}[\log_2 p(X)] = -\sum_{x \in X} p(x) \log_2 p(x)$$



$$-\frac{8}{9} \log_2 \frac{8}{9} - \frac{1}{9} \log_2 \frac{1}{9} \approx \frac{1}{2}$$

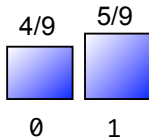
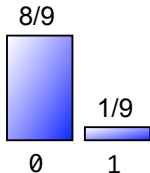


$$-\frac{4}{9} \log_2 \frac{4}{9} - \frac{5}{9} \log_2 \frac{5}{9} \approx 0.99$$

Lec 7: Quantifying Uncertainty

Entropy is a measure of expected “surprise”: How uncertain are we of the value of a draw from this distribution?

$$H(X) = -\mathbb{E}_{X \sim p}[\log_2 p(X)] = -\sum_{x \in X} p(x) \log_2 p(x)$$



$$-\frac{8}{9} \log_2 \frac{8}{9} - \frac{1}{9} \log_2 \frac{1}{9} \approx \frac{1}{2} \quad -\frac{4}{9} \log_2 \frac{4}{9} - \frac{5}{9} \log_2 \frac{5}{9} \approx 0.99$$

- We treat frequency as probability.
- Averages over information content of each observation
- Unit = **bits** (based on the base of logarithm)
- A fair coin flip has 1 bit of entropy

Lec 7: Entropy

- **“High Entropy”**:
 - ▶ Variable has a uniform like distribution
 - ▶ Flat histogram
 - ▶ Values sampled from it are less predictable
- **“Low Entropy”**
 - ▶ Distribution of variable has peaks and valleys
 - ▶ Histogram has lows and highs
 - ▶ Values sampled from it are more predictable
- How much information about one variable do we get by discovering the other variable?

$$IG(Y|X) = H(Y) - H(Y|X)$$

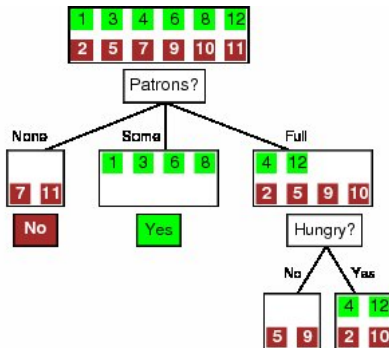
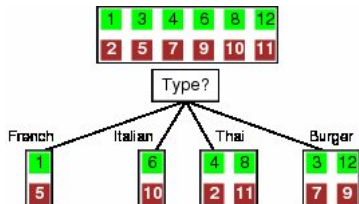
Lec 7: Entropy

- “**High Entropy**”:
 - ▶ Variable has a uniform like distribution
 - ▶ Flat histogram
 - ▶ Values sampled from it are less predictable
- “**Low Entropy**”
 - ▶ Distribution of variable has peaks and valleys
 - ▶ Histogram has lows and highs
 - ▶ Values sampled from it are more predictable
- How much information about one variable do we get by discovering the other variable?

$$IG(Y|X) = H(Y) - H(Y|X)$$

- This is called the **information gain** in Y due to X , or the **mutual information** of Y and X
- If X is completely uninformative about Y : $IG(Y|X) = 0$
- If X is completely informative about Y : $IG(Y|X) = H(Y)$

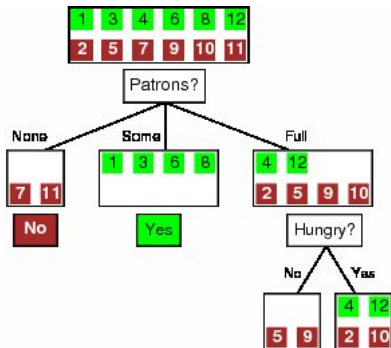
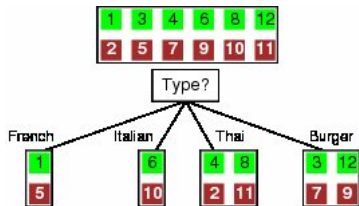
Lec 7: Choosing the best split



$$IG(Y) = H(Y) - H(Y|X)$$

$$IG(\text{type}) = 1 - \left[\frac{2}{12} H(Y|\text{Fr.}) + \frac{2}{12} H(Y|\text{It.}) + \frac{4}{12} H(Y|\text{Thai}) + \frac{4}{12} H(Y|\text{Bur.}) \right] = 0$$

Lec 7: Choosing the best split



$$IG(Y) = H(Y) - H(Y|X)$$

$$IG(\text{type}) = 1 - \left[\frac{2}{12} H(Y|\text{Fr.}) + \frac{2}{12} H(Y|\text{It.}) + \frac{4}{12} H(Y|\text{Thai}) + \frac{4}{12} H(Y|\text{Bur.}) \right] = 0$$

$$IG(\text{Patrons}) = 1 - \left[\frac{2}{12} H(Y|\text{None}) + \frac{4}{12} H(Y|\text{Some}) + \frac{6}{12} H(Y|\text{Full}) \right] \approx 0.541$$

Lec 7: Ensembles : Bagging

- Recall bias-variance trade-off from first half.
- Ensembles combine classifiers to improve performance
- Bagging
 - ▶ Reduces variance (large ensemble can't cause overfitting)
 - ▶ Bias is not changed
 - ▶ Parallel computation
 - ▶ Want to minimize correlation between ensemble elements.
- Random forests are bagged decision trees with an extra randomness: at each split choose a random subset of the attributes.

Lec 8: k-Means

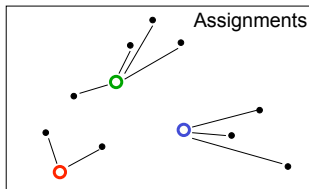
High level overview of algorithm:

- **Initialization:** randomly initialize cluster centers

Lec 8: k-Means

High level overview of algorithm:

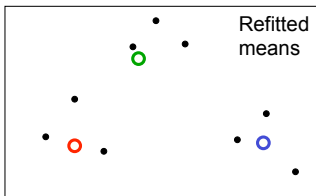
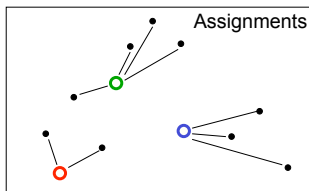
- **Initialization:** randomly initialize cluster centers
- The algorithm iteratively alternates between two steps:
 - ▶ **Assignment step (E-step):** Assign each data point to the closest cluster



Lec 8: k-Means

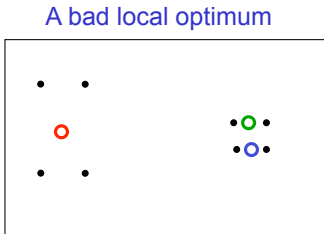
High level overview of algorithm:

- **Initialization:** randomly initialize cluster centers
- The algorithm iteratively alternates between two steps:
 - ▶ **Assignment step (E-step):** Assign each data point to the closest cluster
 - ▶ **Refitting step (M-step):** Move each cluster center to the mean of the data assigned to it



Lec 8: k-Means Algorithm

- The objective is non-convex (so coordinate descent is not guaranteed to converge to the global minimum)
- There is nothing to prevent k-means getting stuck at local minima.
- We could try many random starting points



Lec 8: EM Algorithm

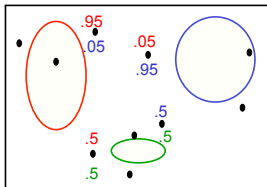
- **Expectation-Maximization alg** alternates between 2 steps:
 1. **E-step**: Compute the posterior probabilities (responsibilities) $\gamma_{nk} = p(z_{nk} = 1 | \mathbf{x}_n)$ given our current model - i.e. how much do we think a cluster is responsible for generating a datapoint.

Lec 8: EM Algorithm

- **Expectation-Maximization alg** alternates between 2 steps:
 1. **E-step:** Compute the posterior probabilities (responsibilities) $\gamma_{nk} = p(z_{nk} = 1 | x_n)$ given our current model - i.e. how much do we think a cluster is responsible for generating a datapoint.
 2. **M-step:** Use the equations derived in lec 8 to update the parameters, assuming γ_{nk} are held fixed- change the parameters of each Gaussian to maximize the probability that it would generate the data it is currently responsible for.

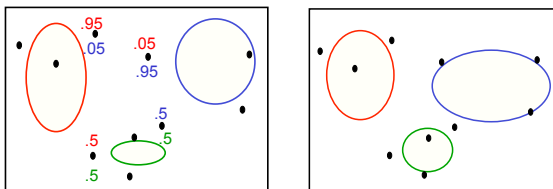
Lec 8: EM Algorithm

- **Expectation-Maximization alg** alternates between 2 steps:
 1. **E-step:** Compute the posterior probabilities (responsibilities) $\gamma_{nk} = p(z_{nk} = 1 | x_n)$ given our current model - i.e. how much do we think a cluster is responsible for generating a datapoint.
 2. **M-step:** Use the equations derived in lec 8 to update the parameters, assuming γ_{nk} are held fixed- change the parameters of each Gaussian to maximize the probability that it would generate the data it is currently responsible for.



Lec 8: EM Algorithm

- **Expectation-Maximization alg** alternates between 2 steps:
 1. **E-step**: Compute the posterior probabilities (responsibilities) $\gamma_{nk} = p(z_{nk} = 1 | x_n)$ given our current model - i.e. how much do we think a cluster is responsible for generating a datapoint.
 2. **M-step**: Use the equations derived in lec 8 to update the parameters, assuming γ_{nk} are held fixed- change the parameters of each Gaussian to maximize the probability that it would generate the data it is currently responsible for.



- When used with general covariance, it can capture the shape of data. K-means with Euclidean distance cannot.

Lec 9: Principal Component Analysis

- Dimensionality reduction: map data to a lower dimensional space
 - ▶ Save computation/memory
 - ▶ Reduce overfitting, achieve better generalization
 - ▶ Visualize in 2 dimensions
- We can approach this problem by:
 - ▶ Minimize the **distortion**: Find vectors \tilde{x} in a subspace that are closest to data points.

$$\min_{\mathbf{U}} \frac{1}{N} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2$$

Lec 9: Principal Component Analysis

- Dimensionality reduction: map data to a lower dimensional space
 - ▶ Save computation/memory
 - ▶ Reduce overfitting, achieve better generalization
 - ▶ Visualize in 2 dimensions
- We can approach this problem by:
 - ▶ Minimize the **distortion**: Find vectors \tilde{x} in a subspace that are closest to data points.

$$\min_{\mathbf{U}} \frac{1}{N} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2$$

- ▶ Maximize the **variance of reconstructions**: Find a direction where data has the most variability.

$$\max_{\mathbf{u}} \frac{1}{N} \sum_i (x_i^T \mathbf{u} - \hat{\mu}^T \mathbf{u})^2$$

- ▶ Both of them lead to PCA.

Lec 9: Principal Component Analysis

- Consider the **empirical covariance matrix**:

$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^\top$$

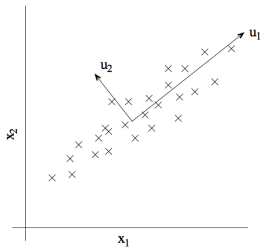
- Recall: Covariance matrices are symmetric and positive semidefinite.

Lec 9: Principal Component Analysis

- Consider the **empirical covariance matrix**:

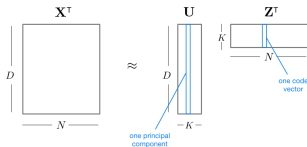
$$\hat{\Sigma} = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})(x_i - \hat{\mu})^T$$

- Recall: Covariance matrices are symmetric and positive semidefinite.
- The optimal PCA subspace is spanned by the top K eigenvectors of $\hat{\Sigma}$.
 - More precisely, choose the first K of any orthonormal eigenbasis for $\hat{\Sigma}$.
- These eigenvectors are called **principal components**, analogous to the principal axes of an ellipse.



Lec 9: PCA as Matrix Factorization

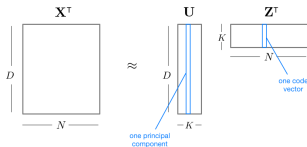
- PCA is approximating $\mathbf{X} \approx \mathbf{Z}\mathbf{U}^T$, or equivalently $\mathbf{X}^T \approx \mathbf{U}\mathbf{Z}^T$.



- Based on the sizes of the matrices, this is a rank- K approximation.

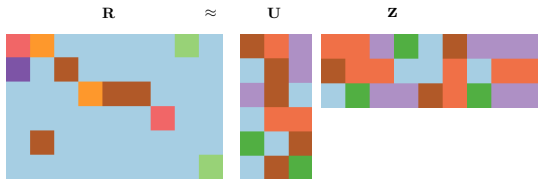
Lec 9: PCA as Matrix Factorization

- PCA is approximating $\mathbf{X} \approx \mathbf{Z}\mathbf{U}^T$, or equivalently $\mathbf{X}^T \approx \mathbf{U}\mathbf{Z}^T$.



- Based on the sizes of the matrices, this is a rank- K approximation.
- Since \mathbf{U} was chosen to minimize reconstruction error, this is the *optimal* rank- K approximation, in terms of error $\|\mathbf{X}^T - \mathbf{U}\mathbf{Z}^T\|_F^2$.

Lec 9: Recommender Systems



- How do we enforce that rating matrix satisfies $R \approx UZ^T$

- ▶ Try

$$\min_{U, Z} \sum_{i, j} (R_{ij} - u_i^T z_j)^2$$

- ▶ Most entries of R are missing! (unseen movies)
- Let $O = \{(n, m) : \text{entry } (n, m) \text{ of matrix } R \text{ is observed}\}$
- Using the squared error loss, a matrix factorization corresponds to solving

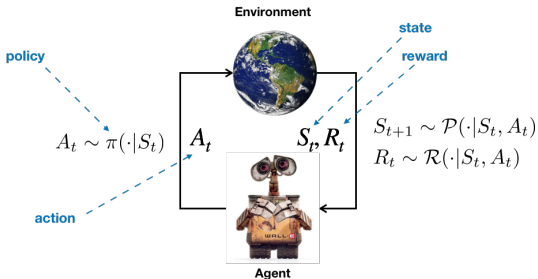
$$\min_{U, Z} \frac{1}{2} \sum_{(n, m) \in O} (R_{nm} - u_n^T z_m)^2$$

- Use alternating least squares, stochastic gradient descent, etc to minimize it.

Lec 10: Reinforcement Learning

A discounted MDP is defined by a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$.

- \mathcal{S} : State space. Discrete or continuous
- \mathcal{A} : Action space. Here we consider finite action space, i.e., $\mathcal{A} = \{a_1, \dots, a_M\}$.
- \mathcal{P} : Transition probability
- \mathcal{R} : Immediate reward distribution
- γ : Discount factor ($0 \leq \gamma \leq 1$)



Lec 10: Reinforcement Learning

The value function is the expected discounted reward if the agent starts from state s , takes action a , and afterwards follows policy π , and satisfies the following recursive relationship:

$$\text{Bellman equations: } Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s, A_0 = a \right]$$

Lec 10: Reinforcement Learning

The value function is the expected discounted reward if the agent starts from state \mathbf{s} , takes action \mathbf{a} , and afterwards follows policy π , and satisfies the following recursive relationship:

$$\text{Bellman equations: } Q^\pi(\mathbf{s}, \mathbf{a}) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = \mathbf{s}, A_0 = \mathbf{a} \right]$$

$$\text{Bellman operator: } = \underbrace{r(\mathbf{s}, \mathbf{a}) + \gamma \int_{\mathcal{S}} Q^\pi(\mathbf{s}', \pi(\mathbf{s}')) \mathcal{P}(\mathbf{s}' | \mathbf{s}, \mathbf{a}) d\mathbf{s}'}_{\triangleq (T^\pi Q^\pi)(\mathbf{s}, \mathbf{a})}$$

Lec 10: Reinforcement Learning

The value function is the expected discounted reward if the agent starts from state s , takes action a , and afterwards follows policy π , and satisfies the following recursive relationship:

$$\text{Bellman equations: } Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s, A_0 = a \right]$$

$$\text{Bellman operator: } = r(s, a) + \underbrace{\gamma \int_{\mathcal{S}} Q^\pi(s', \pi(s')) \mathcal{P}(s' | s, a) ds'}_{\triangleq (T^\pi Q^\pi)(s, a)}$$

$$\text{Bellman Opt Operator: } (T^* Q)(s, a) = r(s, a) + \gamma \int_{\mathcal{S}} \max_{a' \in \mathcal{A}} Q(s', a') \mathcal{P}(s' | s, a) ds'$$

Lec 10: Reinforcement Learning

The value function is the expected discounted reward if the agent starts from state s , takes action a , and afterwards follows policy π , and satisfies the following recursive relationship:

$$\text{Bellman equations: } Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s, A_0 = a \right]$$

$$\text{Bellman operator: } = r(s, a) + \underbrace{\gamma \int_{\mathcal{S}} Q^\pi(s', \pi(s')) \mathcal{P}(s' | s, a) ds'}_{\triangleq (T^\pi Q^\pi)(s, a)}$$

$$\text{Bellman Opt Operator: } (T^* Q)(s, a) = r(s, a) + \gamma \int_{\mathcal{S}} \max_{a' \in \mathcal{A}} Q(s', a') \mathcal{P}(s' | s, a) ds'$$

- Key observation for

$$Q^*(s, a) = \sup_{\pi} Q^\pi(s, a) = \max_{\pi} Q^\pi(s, a):$$

$$Q^* = T^* Q^*$$

Lec 10: Reinforcement Learning

The value function is the expected discounted reward if the agent starts from state s , takes action a , and afterwards follows policy π , and satisfies the following recursive relationship:

$$\text{Bellman equations: } Q^\pi(s, a) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s, A_0 = a \right]$$

$$\text{Bellman operator: } = \underbrace{r(s, a) + \gamma \int_S Q^\pi(s', \pi(s')) \mathcal{P}(s' | s, a) ds'}_{\triangleq (T^\pi Q^\pi)(s, a)}$$

$$\text{Bellman Opt Operator: } (T^* Q)(s, a) = r(s, a) + \gamma \int_S \max_{a' \in \mathcal{A}} Q(s', a') \mathcal{P}(s' | s, a) ds'$$

- Key observation for

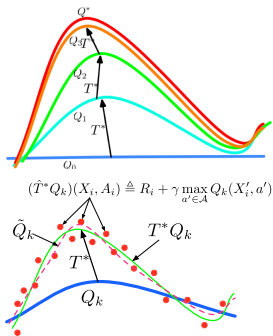
$$Q^*(s, a) = \sup_\pi Q^\pi(s, a) = \max_\pi Q^\pi(s, a):$$

$$Q^* = T^* Q^*$$

- **If we find a Q such that $T^* Q = Q$, then $Q = Q^*$.** We just need to find a fixed point of the operator T^* .

Lec 10: Reinforcement Learning

Three algorithms:



$$(\hat{T}^* Q_k)(X_i, A_i) \triangleq R_i + \gamma \max_{a' \in \mathcal{A}} Q_k(X'_i, a')$$

1. Value iteration: $Q_{k+1} \leftarrow T^* Q_k$

It assumes the **model is known** (distribution $\mathcal{P}(\cdot|s, a)$)

2. Batch RL: **Given the dataset** $\mathcal{D}_N = \{(S_i, A_i, R_i, S'_i)\}_{i=1}^N$ and an action-value function estimate Q_k , we solve a regression problem. We minimize the squared error:

$$Q_{k+1} \leftarrow \underset{Q \in \mathcal{F}}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \left| Q(S_i, A_i) - \left(R_i + \gamma \max_{a' \in \mathcal{A}} Q(S'_i, a') \right) \right|^2$$

3. Q-learning: **Online** update for the action-value function at state S_t :

$$A_t \leftarrow \pi_\epsilon(S; Q) = \begin{cases} \operatorname{argmax}_{a \in \mathcal{A}} Q(S, a) & \text{with probability } 1 - \epsilon \\ \text{Uniformly random action in } \mathcal{A} & \text{with probability } \epsilon \end{cases}$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left[R_t + \gamma \max_{a' \in \mathcal{A}} Q(S_{t+1}, a') - Q(S_t, A_t) \right]$$

STA 414/2104 Closing Remarks

STA 414/2104 Closing Remarks

Continuing with machine learning

- Courses
 - ▶ CSC 412/2506, “Probabilistic Graphical Models”
 - ▶ CSC 421/2516, “Neural Networks and Deep Learning”
 - ▶ CSC 2515, “Machine Learning”
 - ▶ CSC 2532, “Statistical Learning Theory”
 - ▶ Topics courses (varies from year to year): Reinforcement Learning, Algorithmic Fairness, Computer Vision w/ ML, NLP w/ ML, Health w/ ML etc.

STA 414/2104 Closing Remarks

Continuing with machine learning

- Courses
 - ▶ CSC 412/2506, “Probabilistic Graphical Models”
 - ▶ CSC 421/2516, “Neural Networks and Deep Learning”
 - ▶ CSC 2515, “Machine Learning”
 - ▶ CSC 2532, “Statistical Learning Theory”
 - ▶ Topics courses (varies from year to year): Reinforcement Learning, Algorithmic Fairness, Computer Vision w/ ML, NLP w/ ML, Health w/ ML etc.
- Videos from top ML conferences (NeurIPS, ICML, ICLR)

STA 414/2104 Closing Remarks

Continuing with machine learning

- Courses
 - ▶ CSC 412/2506, “Probabilistic Graphical Models”
 - ▶ CSC 421/2516, “Neural Networks and Deep Learning”
 - ▶ CSC 2515, “Machine Learning”
 - ▶ CSC 2532, “Statistical Learning Theory”
 - ▶ Topics courses (varies from year to year): Reinforcement Learning, Algorithmic Fairness, Computer Vision w/ ML, NLP w/ ML, Health w/ ML etc.
- Videos from top ML conferences (NeurIPS, ICML, ICLR)
- Try to reproduce results from papers
 - ▶ If they’ve released code, you can use that as a guide if you get stuck
- Lots of excellent free resources available online!

Thank you!