10-301/601: Introduction to Machine Learning Lecture 15 – Learning Theory (Infinite Case)

Henry Chai & Matt Gormley & Hoda Heidari 10/23/23

Front Matter

- Announcements
 - HW5 released 10/9, due 10/27 (Friday) at 11:59 PM
 - Exam 3 scheduled
 - Tuesday, December 12th from 5:30 PM to 8:30 PM
 - Sign up for peer tutoring! See <u>Piazza</u> for more details

Recall -Theorem 1: Finite, Realizable Case For a *finite* hypothesis set *H* such that c^{*} ∈ *H* (*realizable*) and arbitrary distribution p^{*}, if the number of labelled training data points satisfies

$$N \ge \frac{1}{\epsilon} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ with $\widehat{R}(h) = 0$ have $R(h) \le \epsilon$

Recall -Theorem 1: Finite, Realizable Case For a *finite* hypothesis set *H* such that c^{*} ∈ *H* (*realizable*) and arbitrary distribution p^{*}, if the number of labelled training data points satisfies

$$N = \frac{1}{\epsilon} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ with $\widehat{R}(h) = 0$ have $R(h) \le \epsilon$

• Making the bound tight and solving for ϵ gives...

Statistical Learning Theory Corollary • For a *finite* hypothesis set \mathcal{H} such that $c^* \in \mathcal{H}$ (*realizable*) and arbitrary distribution p^* , given a training dataset S where |S| = N, all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \leq \frac{1}{N} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{1}{\delta}\right) \right)$

with probability at least $1 - \delta$.

Recall -Theorem 2: Finite, Agnostic Case • For a *finite* hypothesis set \mathcal{H} and arbitrary distribution p^* , if the number of labelled training data points satisfies

$$N \ge \frac{1}{2\epsilon^2} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ satisfy $|R(h) - \hat{R}(h)| \le \epsilon$

• Bound is inversely quadratic in ϵ , e.g., halving ϵ means we need four times as many labelled training data points

Statistical Learning Theory Corollary • For a *finite* hypothesis set \mathcal{H} and arbitrary distribution p^* , given a training dataset S where |S| = N, all $h \in \mathcal{H}$ have

$$R(h) \leq \hat{R}(h) + \sqrt{\frac{1}{2N} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right)}$$

with probability at least $1 - \delta$.

What happens when $|\mathcal{H}| = \infty$?

• For a *finite* hypothesis set \mathcal{H} and arbitrary distribution p^* , given a training data set S where |S| = N, all $h \in \mathcal{H}$ have

$$R(h) \le \hat{R}(h) + \sqrt{\frac{1}{2N} \left(\ln(|\mathcal{H}|) + \ln\left(\frac{2}{\delta}\right) \right)}$$

with probability at least $1 - \delta$.

Labellings

- Given some finite set of data points $S = \{x^{(1)}, ..., x^{(N)}\}$ and some hypothesis $h \in \mathcal{H}$, applying h to each point in S results in a <u>labelling</u>
 - [h(x⁽¹⁾), ..., h(x^(N))] is a vector of N +1's and -1's (recall: our discussion of PAC learning assumes binary classification)
- Given $S = \{x^{(1)}, \dots, x^{(N)}\}$, each hypothesis in \mathcal{H}

induces a labelling but not necessarily a unique labelling

• The set of labellings induced by \mathcal{H} on S is

 $\mathcal{H}(S) = \left\{ \left[h(\boldsymbol{x}^{(1)}), \dots, h(\boldsymbol{x}^{(N)}) \right] \mid h \in \mathcal{H} \right\}$

Example: Labellings

 $\mathcal{H} = \{h_1, h_2, h_3\}$



 h_1 h_2



 $\mathcal{H} = \{h_1, h_2, h_3\}$

 $[h_1(\mathbf{x}^{(1)}), h_1(\mathbf{x}^{(2)}), h_1(\mathbf{x}^{(3)}), h_1(\mathbf{x}^{(4)})]$ = (-1, +1, -1, +1)





 $\mathcal{H} = \{h_1, h_2, h_3\}$

 $[h_1(\mathbf{x}^{(1)}), h_1(\mathbf{x}^{(2)}), h_1(\mathbf{x}^{(3)}), h_1(\mathbf{x}^{(4)})]$ = (-1, +1, -1, +1)





 $\mathcal{H} = \{h_1, h_2, h_3\}$

 $\begin{bmatrix} h_1(\mathbf{x}^{(1)}), h_1(\mathbf{x}^{(2)}), h_1(\mathbf{x}^{(3)}), h_1(\mathbf{x}^{(4)}) \end{bmatrix}$ = (+1, +1, -1, -1)





 $|\mathcal{H}(S)| = 2$









 h_1 h_2

VC-Dimension

• $\mathcal{H}(S)$ is the set of all labellings induced by \mathcal{H} on S

- If |S| = N, then $|\mathcal{H}(S)| \le 2^N$
- \mathcal{H} shatters *S* if $|\mathcal{H}(S)| = 2^N$
- The <u>VC-dimension</u> of \mathcal{H} , $VC(\mathcal{H})$, is the size of the largest set *S* that can be shattered by \mathcal{H} .
 - If \mathcal{H} can shatter arbitrarily large finite sets, then $VC(\mathcal{H}) = \infty$
- To prove that $VC(\mathcal{H}) = d$, you need to show
 - 1. \exists some set of d data points that \mathcal{H} can shatter and
 - 2. \nexists a set of d + 1 data points that \mathcal{H} can shatter

• $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

• What is $VC(\mathcal{H})$?

• Can $\mathcal H$ shatter some set of 1 point?





• $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

• What is $VC(\mathcal{H})$?

- Can $\mathcal H$ shatter some set of 1 point?
- Can $\mathcal H$ shatter some set of 2 points?



- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $VC(\mathcal{H})$?
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can \mathcal{H} shatter some set of 2 points?
 - Can \mathcal{H} shatter some set of 3 points?



- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $VC(\mathcal{H})$?
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can \mathcal{H} shatter some set of 2 points?
 - Can \mathcal{H} shatter some set of 3 points?



- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $VC(\mathcal{H})$?
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can \mathcal{H} shatter some set of 2 points?
 - Can \mathcal{H} shatter some set of 3 points?



- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $VC(\mathcal{H})$?
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can \mathcal{H} shatter some set of 2 points?
 - Can \mathcal{H} shatter **some** set of 3 points?



- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $VC(\mathcal{H})$?
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can \mathcal{H} shatter some set of 2 points?
 - Can \mathcal{H} shatter some set of 3 points?



- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $VC(\mathcal{H})$?
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can \mathcal{H} shatter some set of 2 points?
 - Can \mathcal{H} shatter some set of 3 points?



- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $VC(\mathcal{H})$?
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can $\mathcal H$ shatter some set of 2 points?
 - Can \mathcal{H} shatter some set of 3 points?
 - Can $\mathcal H$ shatter some set of 4 points?





S₁ All points on the convex hull

At least one point inside the convex hull

• $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

• What is $VC(\mathcal{H})$?

- Can $\mathcal H$ shatter some set of 1 point?
- Can \mathcal{H} shatter some set of 2 points?
- Can \mathcal{H} shatter some set of 3 points?
- Can $\mathcal H$ shatter some set of 4 points?





S₁ All points on the convex hull

At least one point inside the convex hull

- $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators
- What is $VC(\mathcal{H})$?
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can $\mathcal H$ shatter some set of 2 points?
 - Can \mathcal{H} shatter some set of 3 points?
 - Can $\mathcal H$ shatter some set of 4 points?





S₁ All points on the convex hull

At least one point inside the convex hull

• $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

• What is $VC(\mathcal{H})$?

- Can $\mathcal H$ shatter some set of 1 point?
- Can $\mathcal H$ shatter some set of 2 points?
- Can \mathcal{H} shatter some set of 3 points?
- Can $\mathcal H$ shatter some set of 4 points?





 $|\mathcal{H}(S_1)| = 14$ All points on the convex hull

At least one point inside the convex hull

• $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

• What is $VC(\mathcal{H})$?

- Can $\mathcal H$ shatter some set of 1 point?
- Can \mathcal{H} shatter some set of 2 points?
- Can \mathcal{H} shatter some set of 3 points?
- Can $\mathcal H$ shatter some set of 4 points?





 $|\mathcal{H}(S_1)| = 14$ All points on the convex hull

At least one point inside the convex hull

• $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

• What is $VC(\mathcal{H})$?

- Can $\mathcal H$ shatter some set of 1 point?
- Can \mathcal{H} shatter some set of 2 points?
- Can \mathcal{H} shatter some set of 3 points?
- Can $\mathcal H$ shatter some set of 4 points?





 $|\mathcal{H}(S_1)| = 14$ All points on the convex hull

At least one point inside the convex hull

• $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

• What is $VC(\mathcal{H})$?

- Can $\mathcal H$ shatter some set of 1 point?
- Can \mathcal{H} shatter some set of 2 points?
- Can \mathcal{H} shatter some set of 3 points?
- Can $\mathcal H$ shatter some set of 4 points?





 $|\mathcal{H}(S_1)| = 14$ All points on the convex hull

At least one point inside the convex hull

 $|\mathcal{H}(S_2)| = 14$

• $x \in \mathbb{R}^2$ and $\mathcal{H} =$ all 2-dimensional linear separators

- $VC(\mathcal{H}) = 3$
 - Can $\mathcal H$ shatter some set of 1 point?
 - Can $\mathcal H$ shatter some set of 2 points?
 - Can \mathcal{H} shatter some set of 3 points?
 - Can $\mathcal H$ shatter some set of 4 points?





 $|\mathcal{H}(S_1)| = 14$ All points on the convex hull

At least one point inside the convex hull

 $|\mathcal{H}(S_2)| = 14$

• $x \in \mathbb{R}^d$ and $\mathcal{H} =$ all d-dimensional linear separators

• $VC(\mathcal{H}) = d + 1$

• $x \in \mathbb{R}$ and \mathcal{H} = all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \operatorname{sign}(x - a)$



• $x \in \mathbb{R}$ and \mathcal{H} = all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \operatorname{sign}(x - a)$



• $x \in \mathbb{R}$ and \mathcal{H} = all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \operatorname{sign}(x - a)$


• $x \in \mathbb{R}$ and \mathcal{H} = all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \operatorname{sign}(x - a)$



• $x \in \mathbb{R}$ and \mathcal{H} = all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \operatorname{sign}(x - a)$



• $x \in \mathbb{R}$ and \mathcal{H} = all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \operatorname{sign}(x - a)$



• $x \in \mathbb{R}$ and \mathcal{H} = all 1-dimensional positive rays, i.e., all hypotheses of the form $h(x; a) = \operatorname{sign}(x - a)$



• $VC(\mathcal{H}) = 1$

• $x \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



What is $VC(\mathcal{H})$? A. 0 B. 1 C. 1.5 (TOXIC) D. 2 E. 3

Poll Question 1:

• $x \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



• $x \in \mathbb{R}$ and $\mathcal{H} =$ all 1-dimensional positive intervals



• $VC(\mathcal{H}) = 2$

VC-Dimension: Example

Theorem 3: Vapnik-Chervonenkis (VC)-Bound • Infinite, realizable case: for any hypothesis set \mathcal{H} such that $c^* \in \mathcal{H}$ and arbitrary distribution p^* , if the number of labelled training data points satisfies

$$M = O\left(\frac{1}{\epsilon} \left(VC(\mathcal{H})\log\left(\frac{1}{\epsilon}\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have $R(h) \le \epsilon$

Statistical Learning Theory Corollary 3 • Infinite, realizable case: for any hypothesis set \mathcal{H} such that $c^* \in \mathcal{H}$ and arbitrary distribution p^* , given a training dataset S where |S| = N, all $h \in \mathcal{H}$ with $\hat{R}(h) = 0$ have

$$R(h) \le O\left(\frac{1}{N}\left(VC(\mathcal{H})\log\left(\frac{N}{VC(\mathcal{H})}\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$$

with probability at least $1 - \delta$.

Theorem 4: Vapnik-Chervonenkis (VC)-Bound • Infinite, agnostic case: for any hypothesis set \mathcal{H} and arbitrary distribution p^* , if the number of labelled training data points satisfies

$$N = O\left(\frac{1}{\epsilon^2} \left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)\right)$$

then with probability at least $1 - \delta$, all $h \in \mathcal{H}$ have $|R(h) - \hat{R}(h)| \le \epsilon$

Statistical Learning Theory Corollary 4 • Infinite, agnostic case: for any hypothesis set \mathcal{H} and arbitrary distribution p^* , given a training dataset Swhere |S| = N, all $h \in \mathcal{H}$ have

$$R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{N}\left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)}\right)$$

with probability at least $1 - \delta$.

Approximation Generalization Tradeoff

How well does *h* generalize? $R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{N}\left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)}\right)$ How well does *h* approximate *c**?

Approximation Generalization Tradeoff

Increases as $VC(\mathcal{H})$ increases $R(h) \leq \hat{R}(h) + O\left(\sqrt{\frac{1}{N}\left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)}\right)$ Decreases as $VC(\mathcal{H})$ increases

Can we use this corollary to guide model selection? • Infinite, agnostic case: for any hypothesis set \mathcal{H} and arbitrary distribution p^* , given a training dataset Swhere |S| = N, all $h \in \mathcal{H}$ have

$$R(h) \le \hat{R}(h) + O\left(\sqrt{\frac{1}{N}\left(VC(\mathcal{H}) + \log\left(\frac{1}{\delta}\right)\right)}\right)$$

with probability at least $1 - \delta$.

Learning Theory and Model Selection



Learning Theory and Model Selection



- How can we find this "best tradeoff" for linear separators?
- Use a regularizer! By (effectively) reducing the number of features our model considers, we reduce its VC-dimension.

Learning Theory Learning Objectives You should be able to...

- Identify the properties of a learning setting and assumptions required to ensure low generalization error
- Distinguish true error, train error, test error
- Define PAC and explain what it means to be approximately correct and what occurs with high probability
- Apply sample complexity bounds to real-world machine learning examples
- Theoretically motivate regularization

10-301/601: Introduction to Machine Learning Lecture 15 – Societal Impacts of ML

Henry Chai & Matt Gormley & Hoda Heidari 10/23/23

ML in Societal Applications

8 WAYS MACHINE LEARNING WILL IMPROVE EDUCATION

BY MATTHEW LYNCH / ② JUNE 12, 2018 / 〇 5

Can an Algorithm Tell When Kids Are in Danger

Child protective agencies are haunted when they fail to save kids. Pittsburgh officials believe a new data analysis program is helping them make better judgment calls.

≡ tech

Features Technology Innovation Partner Zone the techies

Home \rangle Features \rangle Emerging tech & innovation Features

Researcher explains how algorithms can create a fairer legal system

Deep learning is being used to predict critical COVID-19 cases

Artificial Intelligence and Accessibility: Examples of a Technology that Serves People with Disabilities





ັຍໂດເມີຍ Your Future Doctor May Not be Human. This Is the Rise of AI in Medicine.

From mental health apps to robot surgeons, artificial intelligence is already changing the practice of medicine.

ROBO RECRUITING

TheUpshot

Can an Algorithm Hire Better Than a Human?

By Claire Cain Miller

20 JAN 2017 Insight Kevin Petrasic | Benjamin Saul

Algorithms and bias: What lenders need to know

The algorithms that power fintech may discriminat can be difficult to anticipate-and financial institut accountable even when alleged discrimination is o unintentional.

Wanted: The 'perfect babysitter.' Must pass AI scan for respect and attitude.

Artificial intelligence is slated to disrupt 4.5 million jobs for African Americans, who have a 10% greater likelihood of automation-based job loss than other workers

Allana Akhtar Oct 7, 2019, 12:57 PM

Misinformation on coronavirus is proving highly contagious

By DAVID KLEPPER July 29, 2020







Email address

(f) 🖾 (r)

ZIP code G

The New Hork Times

I.R.S. Changes Audit Practice That **Discriminated Against Black Taxpayers**

The agency will overhaul how it scrutinizes returns that claim the earned-income tax credit, which is aimed at alleviating poverty.

If you're not a white male, artificial intelligence's use in healthcare could be dangerous By Robert David Hart - July 10, 2017





There's software used across the country to predict future criminals. And it's biased against blacks by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPub

Societal Goals

Foster:

- Productivity and efficiency gains
- Innovation and economic growth
- Due process

. . .

- Consistency
- Traceability
- Making choices & biases evident

Mitigate:

- Violations of human rights
 - Justice, equity, and non-discrimination
 - Privacy and non-surveillance
 - Freedom of communication and expression
 - O Economic freedom
- Negative impact on human flourishing and wellbeing
 - Loss of human sovereignty and control
 - Human cognitive abilities
 - 0

...

Al Incidents on the Rise



Summary statistics Summary visualisations Incidents Articles Evolution of incidents by AI principle V All time total 6264 36345 Privacy & data governance: 256 total 317 1768 Current month's Respect of human rights: 313 300 Number of incidents Transparency & explainability: 325 616 3227 250 Fairness: 163 Peak month 2023-10 2023-10 200 Robustness & digital security: 341 616 3227 Reskill or upskill: 30amount 150 Accountability: 165 nge (month-over-month) 23.2 51.22 100 Human wellbeing: 17 Performance: 94 13.01 13.87 50 Safety: 118% change (year over year) 961.58 690.9 0 2023-07 2015-07 2017-01 2017-07 2020-07 2021-07 2022-01 2022-07 2023-01 2014-01 2015-01 2016-01 2016-07 2018-01 2018-07 2019-07 2020-01 2021-01 2014-07 2019-01 *Note: Percent change is calculated based on preceding full months (i.e. the current month is excluded). Date

Principles

- Fairness
- Accountability
- Transparency
- Safety and reliability
- Privacy
- ...

/ AI Ethics Guidelines Global Inventory

AlgorithmWatch's inventory of principles, voluntary commitments and frameworks for an ethical use of algorithms and AI (work in progress)...





Safe and Effective
Systems



Algorithmic Discrimination Protections



Data Privacy



Notice and Explanation



Human Alternatives, Consideration, and Fallback



Beyond Principles

Concerns around **impact**:

- Economic (IP, Antitrust, labor market effects)
- Sustainability and environmental
- Eroding democratic values
 - misinformation and disinformation

Concerns around the **process**:

- Human sovereignty, autonomy, agency, self-determination
 - Participation
 - Recourse / appeal
 - Mental health
- ...

Unfairness and Discrimination

Amazon scraps secret AI recruiting tool that showed bias against women

Jeffrey Dastin

8 MIN READ

SAN FRANCISCO (Reuters) - Amazon.com Inc's (AMZN.O) machine-learning specialists uncovered a big problem: their new recruiting engine did not like women.



(Outcome) Unfairness

Formal Principle of Distributive Justice:

"Equals should be treated equally, and unequals unequally, in proportion to relevant similarities and differences." [Aristotle, ..., Feinberg'1973]

Working Definition of Outcome Unfairness:

Disparate or unequal allocation of harm/benefit across socially salient, but morally irrelevant groups of people.

Mathematical Notions of Fairness

- **Group** notions
 - Statistical parity
 - Equality of accuracy
 - Equality of false positive/false negative rates
 - Equality of positive/negative predictive value
- Individual notions
 - Treat similar individuals similarly.
- **Counterfactual** notions

Statistical/Demographic Parity

• Equal **selection rate** across different groups:

$$P[Y^{=}1|S = S_1] = P[Y^{=}1|S = S_2]$$

• Equal Employment Opportunity Commission:

"A selection rate for any race, sex, or ethnic group which is less than four-fifths (or 80%) of the rate for the group with the highest rate will generally be regarded by the Federal enforcement agencies as evidence of [discrimination]."

Equality of Accuracy

• Equality of the prediction accuracy (L) across groups:

 $E[L(y^{,} y) | S = s_1] = E[L(y^{,} y) | S = s_2]$

• **Example:** Gender shades (Buolamwini et al.'18)



Equality of FPR/FNR

• Equality of the **False Positive Rate (FPR)** across groups:

P[Y⁼¹ | Y =0, S = s₁]=P[Y⁼¹ | Y =0, S = s₂]

• Equality of the **False Negative Rate (FNR)** across groups:

 $P[Y^{=0}|Y=1, S=s_1]=P[Y^{=0}|Y=1, S=s_2]$

• Equality of **Odds**: equal FNR and FPR simultaneously



Equality of PPV/NPV

• Equality of the **Positive Predictive Value (PPV)**

 $P[Y = 1 | Y^{=1}, S = s_1] = P[Y = 1 | Y^{=1}, S = s_2]$

• Equality of the **Negative Predictive Value (NPV)**

 $P[Y=0|Y^{=}0, S=s_1]=P[Y=0|Y^{=}0, S=s_2]$

• **Predictive Value Parity (PVP):** equal PPV and NPV simultaneously



Common Pros and Cons

- Ignores possible correlation between Y and S.
- Allows for trading off different types of error.
- Allows laziness.
- Doesn't consider practical considerations.
 - e.g., High accuracy difficult to attain for small groups
- ...

Summary of Fairness Notions w. Confusion Matrix

For each group s, form:

	$\hat{Y}=0$	$\hat{Y}=1$
Y=0	a (true negative)	b (false positive)
Y=1	c (false negative)	d (true positive)

across all s.

Individual vs. Group Fairness

- Treating people as individuals, regardless of their group membership.
- Disparate Treatment:

"Similarly situated individuals must be treated similarly."

• Similarity must be defined *with respect to the task at hand.*

Example: movie casting vs. employment decisions in tech sector

Formalizing Individual Fairness

(Dwork et al. 2012):

- d(**x**_i,**x**_j): a metric defining distance between two individuals
- D: a measure of distance between distributions
- A randomized classifier h mapping \mathbf{x} to $\Delta_h(\mathbf{x})$ satisfies the (D, d)-Lipschitz property if $\forall \mathbf{x}_i, \mathbf{x}_j$,

 $D(\Delta_h(\mathbf{x}_i), \Delta_h(\mathbf{x}_j)) \le d(\mathbf{x}_i, \mathbf{x}_j).$

Several problems with the Formulation

- Does not treat **dis**similar individuals **differently**.
- How should we pick d and D?
- Applicable to probabilistic models, only.
- Computationally expensive (O(n²) pairwise constraints)
- ...