# Lecture 13 Distribution Shift

IEMS 402 Statistical Learning

Northwestern

# References

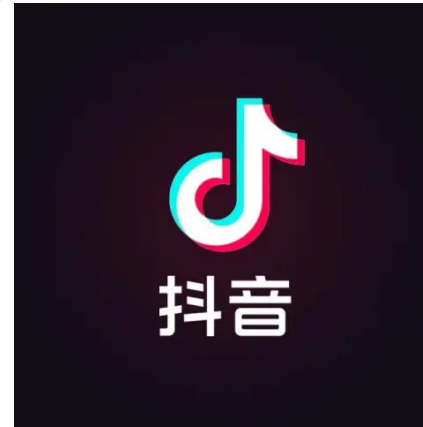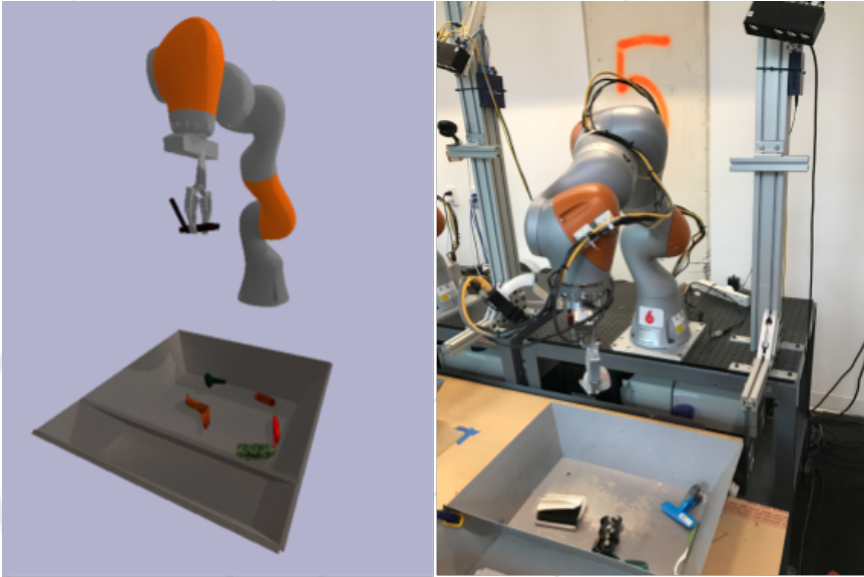https://hsnamkoong.github.io/assets/html/b9145/index.html

# Distribution Shift

# Reconsider the ML Theory…

$$\left| E_{\widehat{\mathbb{P}}} f - E_{\mathbb{P}} f \right| \leq \dots$$

$\widehat{\mathbb{P}}$ is i.i.d. from $\mathbb{P}$.

What if the test distribution $\mathbb{P}_{test}$ is different ??

# However…

# Elephant or Cat

# Shortcut learning

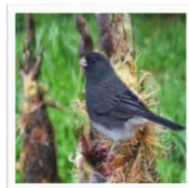| | | | | |
|---|---|---|---|---|
| |  |  |  | **Article:** Super Bowl 50<br><br>**Paragraph:** "*Peython Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager.* Quarterback Jeff Dean had a jersey number 37 in Champ Bowl XXXIV.*"<br><br>**Question:** "*What is the name of the quarterback who was 38 in Super Bowl XXXIII?*"<br><br>**Original Prediction:** John Elway<br>**Prediction under adversary:** Jeff Dean |
| **Task for DNN** | Caption image | Recognise object | Recognise pneumonia | Answer question |
| **Problem** | Describes green hillside as grazing sheep | Hallucinates teapot if certain patterns are present | Fails on scans from new hospitals | Changes answer if irrelevant information is added |
| **Shortcut** | Uses background to recognise primary object | Uses features irrecognisable to humans | Looks at hospital token, not lung | Only looks at last sentence and ignores context |

# spurious correlation
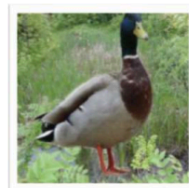


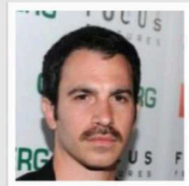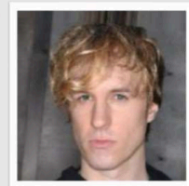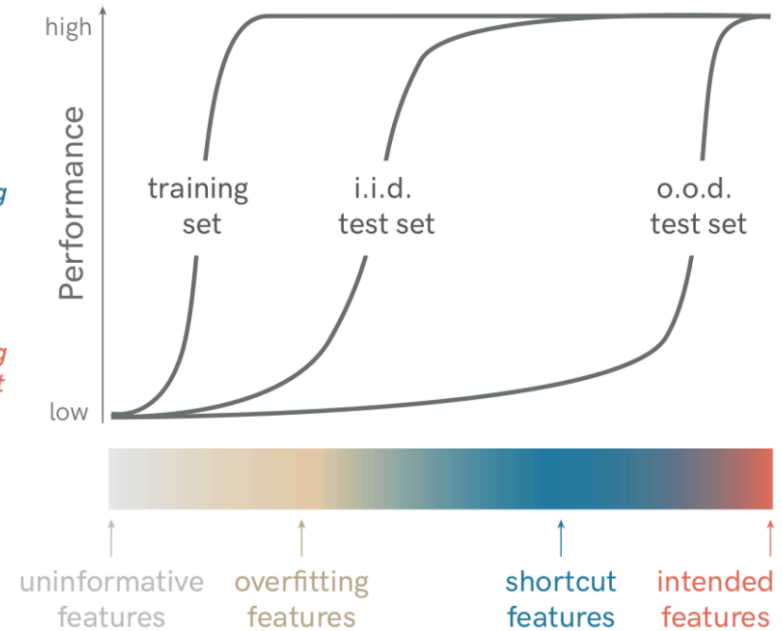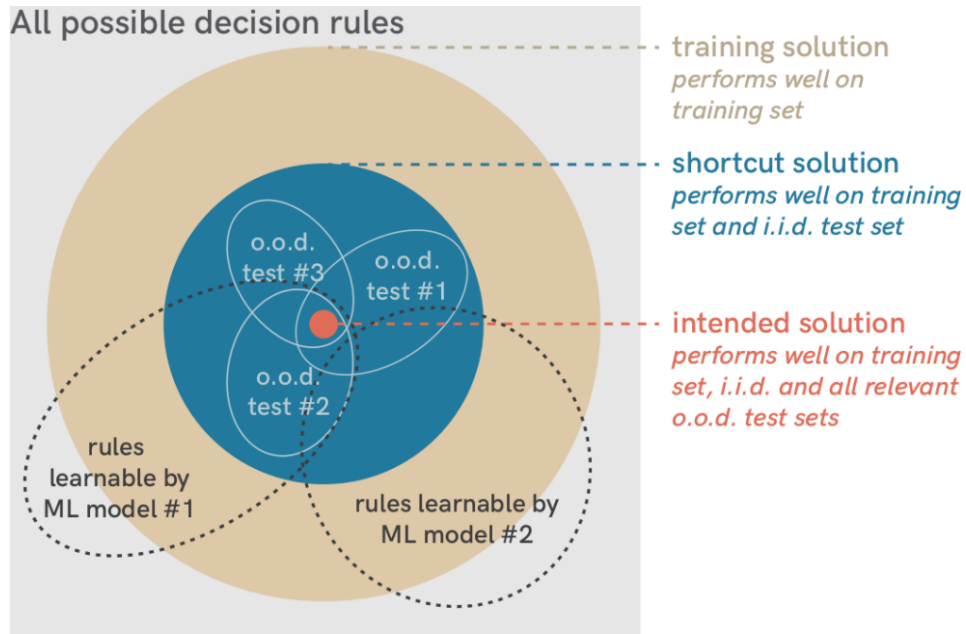|  | Common training examples | | Test examples |
|---|---|---|---|
| **Waterbirds** | y: waterbird<br>a: water<br>background | y: landbird<br>a: land<br>background | y: waterbird<br>a: land<br>background |
| **CelebA** | y: blond hair<br>a: female | y: dark hair<br>a: male | y: blond hair<br>a: male |
| **MultiNLI** | y: contradiction<br>a: has negation<br>(P) The economy could be still better.<br>(H) The economy has never been better. | y: entailment<br>a: no negation<br>(P) Read for Slate's take on Jackson's findings.<br>(H) Slate had an opinion on Jackson's findings. | y: entailment<br>a: has negation<br>(P) There was silence for a moment.<br>(H) There was a short period of time where no one spoke. |

# From i.i.d to o.o.d

# Importance Weighting

# Importance Weighting

How do we deal with covariate / label shifts?

**What we have**

$$E_{p_{train}}[\ell(z; \theta)]$$

**What we want**

$$E_{p_{test}}[\ell(z; \theta)]$$

*Weighted loss.*

Most basic approach: reweight the loss

$$E_{p_{train}}\left[\frac{p_{test}(z)}{p_{train}(z)}\ell(z; \theta)\right] = E_{p_{test}}[\ell(z; \theta)]$$

Weighted loss over the
training distribution

(also possible: resample the dataset)

*reweighting   data 1      data 2*

*① $\sum$ loss (data1) + loss (data2)*

*resample*

*② data1 data2 ... data2 data2*

*in expertation, they are same
but there variant are diffent.*

# Importance weighting

*I don't know $\frac{P_{test}}{P_{train}}$*

An alternative algorithm: use a classifier that separates $p_{train}$ and $p_{test}$

1. Estimate a classifier $f(x) \approx \dfrac{p_{train}(x)}{p_{test(x)} + p_{train}(x)}$

*( collect another dataset.*
*( X_{train}, 1 )    ( X_{test}, 0 ) )*
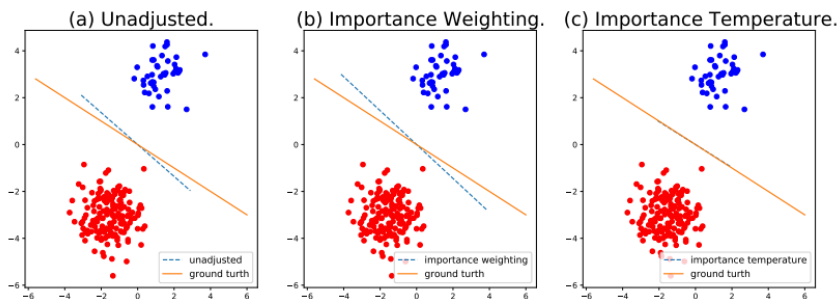
2. Reweight by $h(x) = \dfrac{1}{f(x)} - 1$

3. Fit a model by minimizing the loss $h(x)\ell(x, y; \theta)$

# Not Working for Over-parameterized Model



(a) Linear Model for Separable Data

(b) Multilayer Perceptron with two hidden layers of size 200

Byrd J, Lipton Z. What is the effect of importance weighting in deep learning? International conference on machine learning. PMLR, 2019: 872-881.

# IPM

# Background material: integral probability measures

$$E_{P_{test}} \ell - E_{P_{train}} \ell$$

To state this clearly, we need to first go into some background.

**Definition (IPM):**

For two probability distributions $p$ and $q$, the integral probability metric (IPM) for a family of functions $\mathcal{F}$ is defined as

⟹ uniform

$$d_{\mathcal{F}}(p, q) = \sup_{f \in \mathcal{F}} |E_p[f(x)] - E_q[f(x)]|$$

distance between distribution $P_{train}$ / $P_{test}$

**Intuition:** $\mathcal{F}$ are 'test functions' that can distinguish $p$ and $q$

If two have the same function value for all $\mathcal{F}$, then they are similar

# IPM and distribution shift

What we want       What we have       Domain distance

$$E_{p_{test}}[\ell(x, y, \theta)] = E_{p_{train}}[\ell(x, y, \theta)] + \Delta$$

**From the trivial restatement**

$$\Delta = E_{p_{test}}[\ell(x, y, \theta)] - E_{p_{train}}[\ell(x, y, \theta)]$$

**This looks like an IPM!** (if $\ell(x, y, \theta) \in \mathcal{F}$ for all $\theta$)

$$\Delta \leq \sup_{f \in \mathcal{F}} E_{p_{test}}[f(x, y)] - E_{p_{train}}[f(x, y)] = d_{\mathcal{F}}(p_{train}, p_{test})$$

**Takeaway:** IPMs bound excess error under transfer

$$E_{p_{test}} \ell \leq E_{p_{train}} \ell + d_{\mathcal{F}}(p_{train}, p_{test})$$

# Example: L1 distance

$$d_F(P, \theta) = \sup_{\{f | -1 \leq f \leq 1\}} \left| \sum f(x)(P(x) - \theta(x)) \right|$$
$$= \sum_x |P(x) - \theta(x)|$$

We can now bound test performance in terms of IPMs

$$F := \left\{ f \mid -1 \leq f(x) \leq 1, \forall x \right\}, \quad d_F(P, \theta) = \sum_x |P(x) - \theta(x)|$$

For $0 \leq \ell(x, y, \theta) \leq 1$ and under covariate shift,

$$E_{p_{test}}[\ell(x, y, \theta)] \leq E_{p_{train}}[\ell(x, y, \theta)] + \left\| p_{train}(x) - p_{test}(x) \right\|_1$$



$$\left\| p_{train}(x) - p_{test}(x) \right\|$$

| | **Reweighting** | **IPM** |
|---|---|---|
| **Goals** | Correct train-test mismatch | Estimate train-test mismatch |
| **Assumptions** | Overlap | Boundedness |
| **Training** | Weighted/modified loss | No change |
| **Costs** | More samples (variance) | Inaccurate models (bias) |

Curse of dimensionality (next lecture)

*Handwritten annotations:*

can also correct

$\dfrac{P_{test}}{P_{train}} \Rightarrow$ if $P_{train}(x) = 0$ then $P_{test}(x) = 0$

$d_F(P, \tilde{P})$

you may use $E_{P_{test}} \ell + d_F(P, \tilde{P})$ as loss

# Defining HΔH (disagreement)

For a hypothesis class $\mathcal{H}$, the HΔH set is defined as the symmetric difference



**Definition (HΔH):**

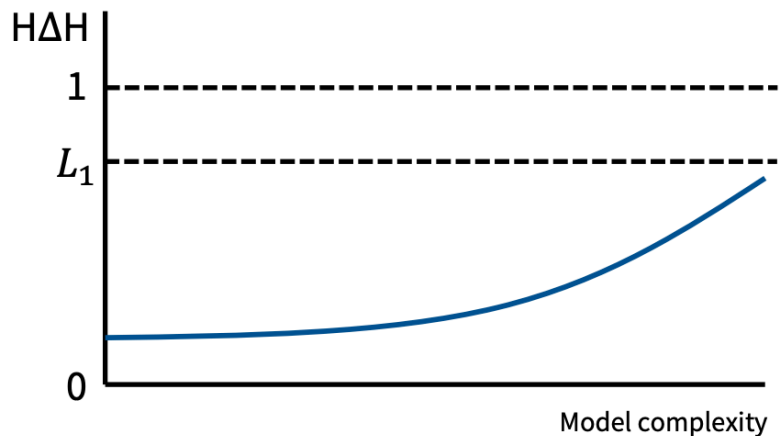For a hypothesis class $\mathcal{H}$, the symmetric difference set $H\Delta H$ is defined as

$$\text{H}\Delta\text{H} \coloneqq \{g\colon g(x) = \text{XOR}\big(h(x), h'(x)\big) \text{ and } h, h' \in \mathcal{H}\}$$

# Dependency on Hypothesis Space

For a hypothesis class $\mathcal{H}$, the HΔH-divergence is

$$d_{H\Delta H}(p_{train}, p_{test}) = 2 \sup_{g \in H\Delta H} \left| E_{p_{train}}[g(x)] - E_{p_{test}}[g(x)] \right|$$

F

**HΔH:** $\frac{1}{2} d_{H\Delta H}(p_{train}, p_{test})$

$d_{H\Delta H}$ is upper bounded by the $L_1$ distance

$d_{H\Delta H}$ increases monotonically with model complexity. If $H \subset H'$,
$$d_{H\Delta H} \leq d_{H'\Delta H'}$$

HΔH

1

$L_1$

0

Model complexity

# Another trade-off

Let's walk through the main bound.

$$E_{p_{test}}[\ell(x, y, h)]$$
$$\leq E_{p_{train}}[\ell(x, y, h)] + \frac{1}{2} d_{H\Delta H}(p_{train}, p_{test}) + \lambda$$

Different answer on two domains

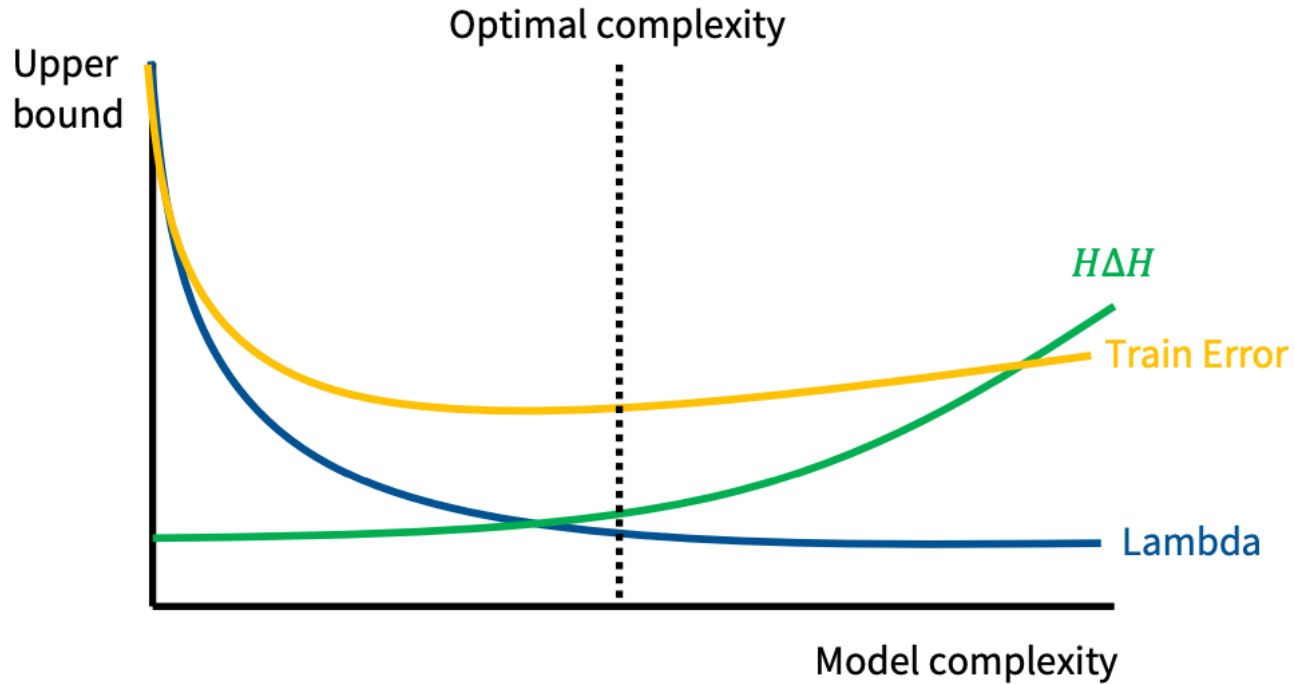same answer but Both are wrong

Training domain error

Domain distinguishability

Minimal error of a classifier on both domains

$$\lambda = \inf_{h \in \mathcal{H}} p_{train}(y \neq h(x)) + p_{test}(y \neq h(x))$$

**HΔH claim:** Low training domain error + low $H\Delta H$ divergence + rich $\mathcal{H}$
= good generalization to target domain

# Another tradeoff

# Distributionally Robust Optimization

# F-divergence

$f(1) = 0$

$\nearrow$ $f$ is a convex function

$f$- divergence.

$$D_f \left( Q \parallel P \right) := \int f\left( \frac{dQ}{dP} \right) dP$$

$f(t) = t \log t,$ then $\quad$ KL- divergence

$f(t) = |t-1| \quad \rightarrow \quad$ 4 distance (Total Variation)

$f(t) = (t-1)^2 \qquad\qquad \chi^2$ divergen

Northwestern

24

# Distributionally Robust Optimization

$$dist(P_{train}, P_{test}) \leq \rho$$

| | |
|---|---|
| **Empirical Risk Minimization** | $\min\limits_{\theta \in \Theta} \mathbb{E}_{Z \sim P_{train}}[\ell(\theta; Z)]$ |

$$\mathbb{E}_{Z \sim P_{test}}[\ell(\theta; Z)]$$

| | |
|---|---|
| **DRO** | $\min\limits_{\theta \in \Theta} \sup\limits_{Q \in \mathcal{P}} \mathbb{E}_{Z \sim Q}[\ell(\theta; Z)]$ |

$$\mathcal{P} = \{Q : Dist(Q, P_{train}) \leq \rho\}$$

$\mathcal{P}$

*distance between distributions*

$P_{train}$

$P_{test}$

Instead of minimizing loss over training distribution, minimize loss over distributions *near* it

# Generalization of DRO

automatically can be built.

$$\sup_{d(Q,P) \leq \rho} \mathbb{E}_{z \sim Q}[\ell(\theta; z)] \geq \mathbb{E}_{z \sim P_{test}}[\ell(\theta; z)]$$
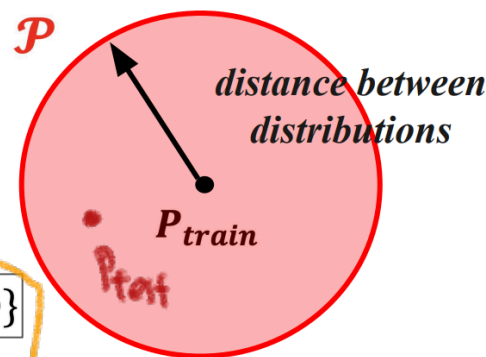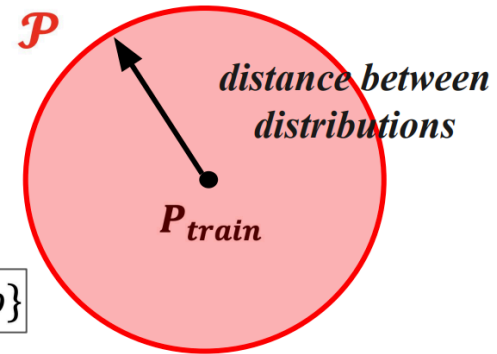
**Empirical Risk Minimization**

$$\min_{\theta \in \Theta} \mathbb{E}_{Z \sim P_{train}}[\ell(\theta; Z)]$$

**DRO**

$$\min_{\theta \in \Theta} \sup_{Q \in \mathcal{P}} \mathbb{E}_{Z \sim Q}[\ell(\theta; Z)]$$

$$\mathcal{P} = \{Q : Dist(Q, P_{train}) \leq \rho\}$$

$\mathcal{P}$

*distance between distributions*

$P_{train}$

Instead of minimizing loss over training distribution, minimize loss over distributions *near* it

# Duality of DRO

$$R_f(\theta; P) := \sup_{D_f(Q\|P) \le \rho} \mathbb{E}_Q[\ell(\theta, z)]$$

$$R_f(\theta; P) = \inf_{\lambda \ge 0, \eta \in \mathbb{R}} \left\{ \lambda \mathbb{E}_P\left[ f^*\left( \frac{\ell(\theta; Z) - \eta}{\lambda} \right) \right] + \lambda\rho + \eta \right\}$$

$$f^*(s) := \sup_t \{st - f(t)\}.$$

$$= \sup_{L \ge 0} \inf_{\lambda \ge 0, \eta \in \mathbb{R}} \left\{ \mathbb{E}_P[L(Z)\ell(\theta; Z)] + \lambda(\rho - \mathbb{E}_P[f(L(Z))] - \eta(\mathbb{E}_P[L(Z)] - 1)) \right\}$$

$$P-D+(Q\|P) \ge 0$$

$$Q \text{ is a distribution.}$$

$$L(2) = \frac{dQ}{dP} \cdot \quad \mathbb{E}_P \frac{dQ}{dP} = \mathbb{E}_Q 1 = 1$$

→ Next step: Solve the sup over L

# Duality of DRO

$$R_f(\theta; P) = \inf_{\lambda \geq 0, \eta \in \mathbb{R}} \left\{ \lambda \mathbb{E}_P \left[ f^* \left( \frac{\ell(\theta; Z) - \eta}{\lambda} \right) \right] + \lambda \rho + \eta \right\}$$

$$f^*(s) := \sup_t \{ st - f(t) \}.$$

$$= \sup_{L \geq 0} \inf_{\lambda \geq 0, \eta \in \mathbb{R}} \left\{ \mathbb{E}_P[L(Z)\ell(\theta; Z)] + \lambda(\rho - \mathbb{E}_P[f(L(Z))] - \eta(\mathbb{E}_P[L(Z)] - 1) \right\}$$

$$= \inf_{\lambda \geq 0, \eta \in \mathbb{R}} \sup_{L \geq 0} \left\{ \lambda \mathbb{E}_P \left[ \frac{L(Z)(\ell(\theta; Z) - \eta)}{\lambda} - f(L(Z)) \right] \right\} + \lambda \rho + \eta.$$

$$= \mathbb{E}_P \left[ f^* \left( \frac{\ell(\theta; Z) - \eta}{\lambda} \right) \right].$$

→ reference

↳ The loss function : $f^*$

$f^*(\text{relativa b})$

Duality of f-divergence DRO, is changing loss function ∈.

# Variance Regularization

$\chi^2$ divergence $\cdot$ $\qquad f = (t-1)^2$

$$\inf_{\substack{\lambda \geq 0 \\ \eta \in \mathbb{R}}} \mathbb{E}_{\mathbb{P}} \left( \frac{\ell(\theta; z) - \eta}{\lambda} \right)^2 + \ell(\theta; z).$$

$\downarrow$

$\eta = \mathbb{E}_{\mathbb{P}} \ell$

$$\lambda \, \mathrm{Var}(\ell) + \mathbb{E}_{\mathbb{P}} \ell$$
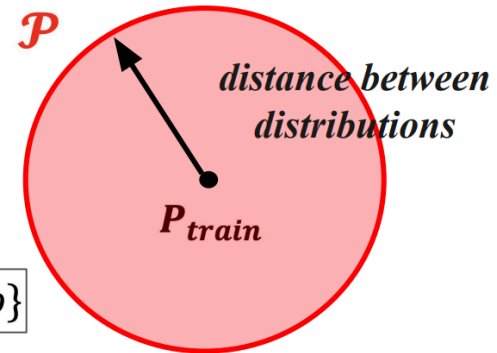
# Generalization of DRO

**Empirical Risk Minimization**

$$\min_{\theta \in \Theta} \mathbb{E}_{Z \sim P_{train}}[\ell(\theta; Z)]$$

**DRO**

$$\min_{\theta \in \Theta} \sup_{Q \in \mathcal{P}} \mathbb{E}_{Z \sim Q}[\ell(\theta; Z)]$$
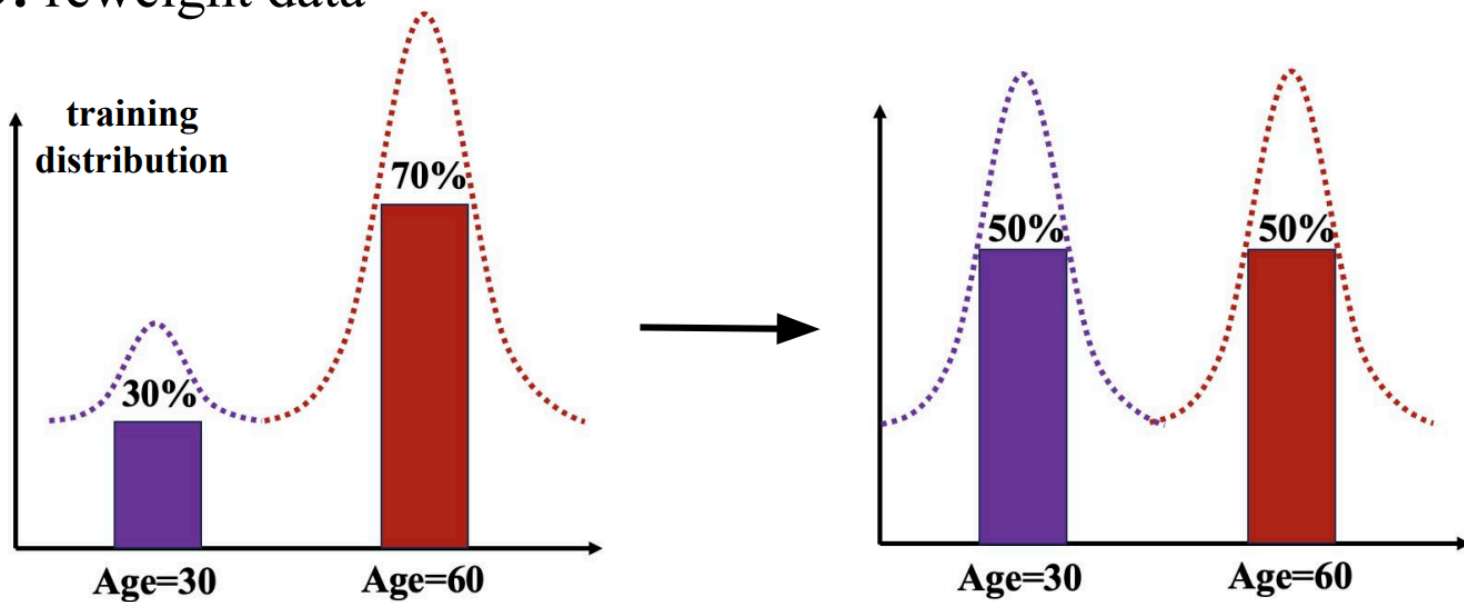
$$\mathcal{P} = \{Q : Dist(Q, P_{train}) \leq \rho\}$$

$\mathcal{P}$

*distance between distributions*

$P_{train}$

Instead of minimizing loss over training distribution, minimize loss over distributions **near** it

# Is DRO Working?

# F-divergence DRO only reweighting



*f*-**DRO:** reweight data

# spurious correlation

weight more on the data have higher loss
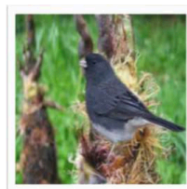
weight more on the data opposing in the first.

Weights more on rare data!

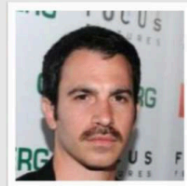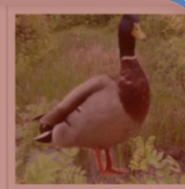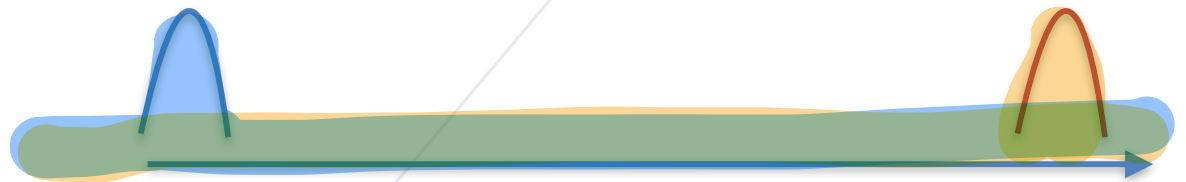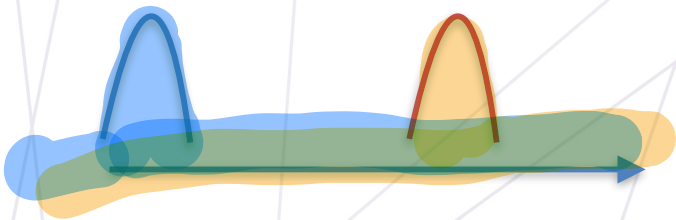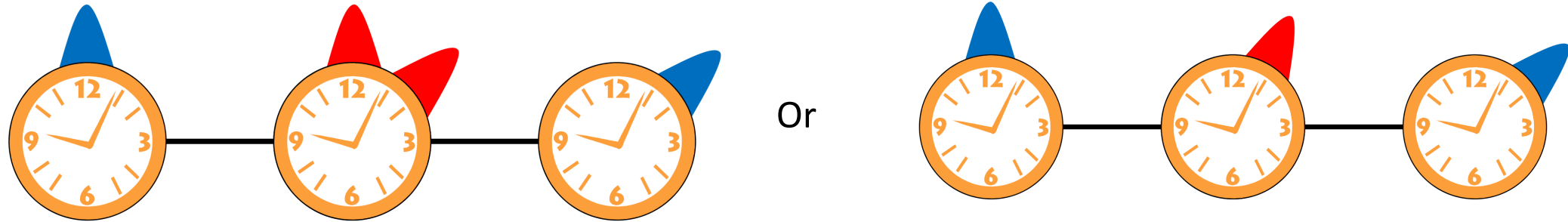|  | **Common training examples** | | **Test examples** |
|---|---|---|---|
| **Waterbirds** | y: waterbird<br>a: water background | y: landbird<br>a: land background | y: waterbird<br>a: land background |
| **CelebA** | y: blond hair<br>a: female | y: dark hair<br>a: male | y: blond hair<br>a: male |
| **MultiNLI** | y: contradiction<br>a: has negation<br>(P) The economy could be still better.<br>(H) The economy has never been better. | y: entailment<br>a: no negation<br>(P) Read for Slate's take on Jackson's findings.<br>(H) Slate had an opinion on Jackson's findings. | y: entailment<br>a: has negation<br>(P) There was silence for a moment.<br>(H) There was a short period of time where no one spoke. |

# What's wrong about f-divergence

$$D_f(P \| Q) = \mathbb{E}_{\sim P} \, f\left(\frac{dQ}{dP}\right)$$

# What's wrong about f-divergence



Or

# Over-parameterization?

| | | | Average Accuracy | | Worst-Group Accuracy | |
|---|---|---|---|---|---|---|
| | | | ERM | DRO | ERM | DRO |
| **Standard Regularization** | Waterbirds | Train | 100.0 | 100.0 | 100.0 | 100.0 |
| | | Test | 97.3 | 97.4 | 60.0 | 76.9 |
| | CelebA | Train | 100.0 | 100.0 | 99.9 | 100.0 |
| | | Test | 94.8 | 94.7 | 41.1 | 41.1 |
| | MultiNLI | Train | 99.9 | 99.3 | 99.9 | 99.0 |
| | | Test | 82.5 | 82.0 | 65.7 | 66.4 |
| **Strong $\ell_2$ Penalty** | Waterbirds | Train | 97.6 | 99.1 | 35.7 | 97.5 |
| | | Test | 95.7 | 96.6 | 21.3 | 84.6 |
| | CelebA | Train | 95.7 | 95.0 | 40.4 | 93.4 |
| | | Test | 95.8 | 93.5 | 37.8 | 86.7 |

*Handwritten annotations: "No improvement" (top right); "model class is m-fl." (lower left); arrows pointing from ERM worst-group to DRO worst-group columns.*

# Adversarial Learning

# adversarial training



"pig"  + 0.005 x  = "airliner"

# How to find Adversarial Examples?



$$+ .007 \times$$

$$\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

$$=$$

$$\boldsymbol{x} + \epsilon \text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

$\boldsymbol{x}$

"panda"
57.7% confidence

"nematode"
8.2% confidence

"gibbon"
99.3 % confidence

Optimization that maximize the loss

# Adversarial Training

$$\min_{\theta} \rho(\theta), \quad \text{where} \quad \rho(\theta) = \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\delta \in \mathcal{S}} L(\theta, x + \delta, y) \right].$$

https://arxiv.org/pdf/1706.06083

# Adversarial Training Can Hurt Generalization

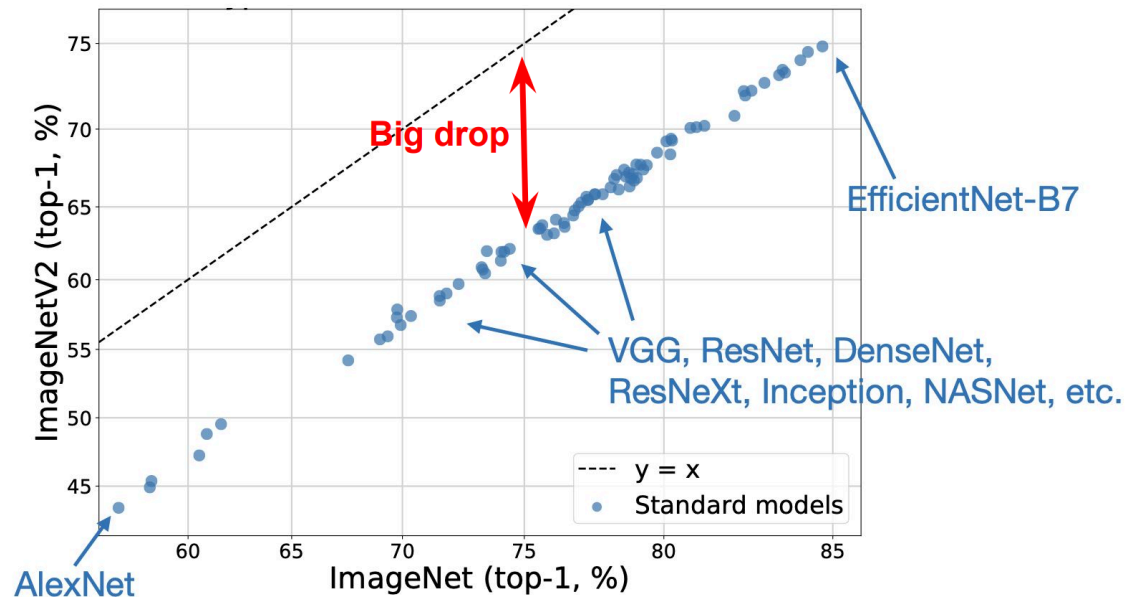|  | Standard training | Adversarial training |
|---|---|---|
| Robust test | 3.5% | 45.8% |
| Robust train | - | 100% |
| Standard test | 95.2% | 87.3% |
| Standard train | 100% | 100% |

# Real World?

Lots of progress on ImageNet over the past 10 years, but models are still not robust.

Evaluation: **new test sets**



**ImageNetV2**

[Recht, Roelofs, Schmidt, Shankar '19]

**ObjectNet**

[Barbu, Mayo, Alverio, Luo, Wang, Gutfreund, Tenenbaum, Katz '19]

**ImageNet-Sketch**

[Wang, Ge, Lipton, Xing '19]

**ImageNet-R**

[Hendrycks, Basart, Mu, Kadavath, Wang, Dorundo, Desai, Zhu, Parajuli, Guo, Song, Steinhardt, Gilmer '20]

# Agree on the line!

Recht B, Roelofs R, Schmidt L, et al. Do imagenet classifiers generalize to imagenet?[C]//
International conference on machine learning. PMLR, 2019: 5389-5400.



[Taori, Dave, Shankar, Carlini, Recht, Schmidt '20]

# Why?