

Training ≠ Testing

- **Generalization performance**
(the empirical distribution lies)
- **Covariate shift**
(the covariate distribution lies)
- **Logistic regression**
(tools to fix shift)
- **Covariate shift correction**
- **Label shift**
(the label distribution lies)
- **Nonstationary Environments**

$$p_{\text{emp}}(x, y) \neq p(x, y)$$

$$p(x) \neq q(x)$$

$$\log(1 + \exp(-yf(x)))$$

$$\frac{1}{2} (p(x)\delta(1, y) + q(x)\delta(-1, y))$$

$$p(y) \neq q(y)$$



Generalization performance

Generalization performance



Generalization performance



Generalization performance



Generalization performance

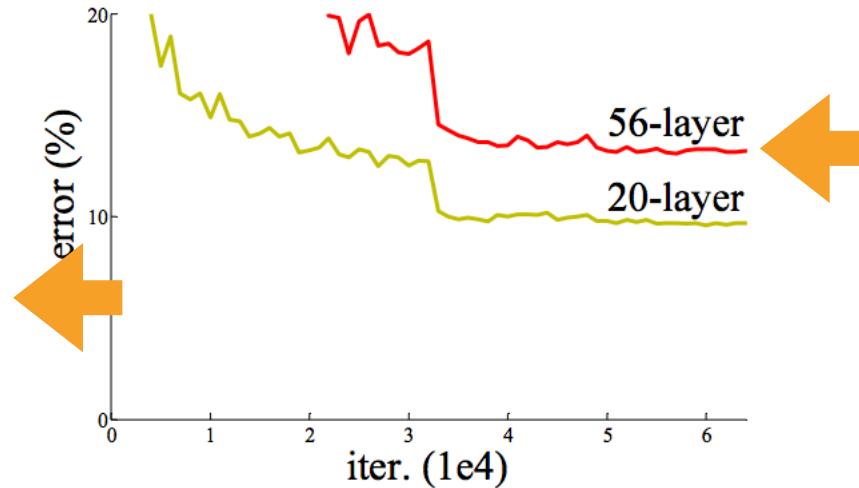
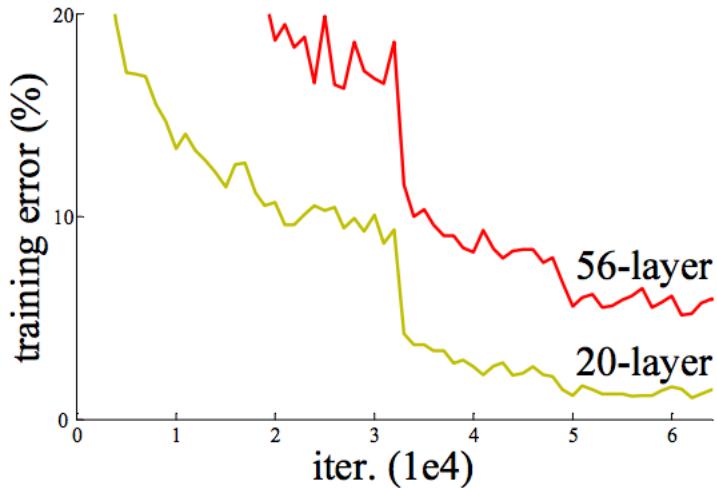


Generalization performance



Only cats and dogs?

- Images, too (e.g. He et al., 2015, ResNet paper)



- Alexa
(‘Please turn off the coffee machine’ vs. ‘coffee machine off’)

Why?

- Data Distribution $p(x,y)$
- Dataset drawn from $p(x,y)$
- Training minimizes empirical risk (plus regularization)

$$\underset{w}{\text{minimize}} \frac{1}{m} \sum_{I=1}^m l(f(x_i, w), y_i)$$

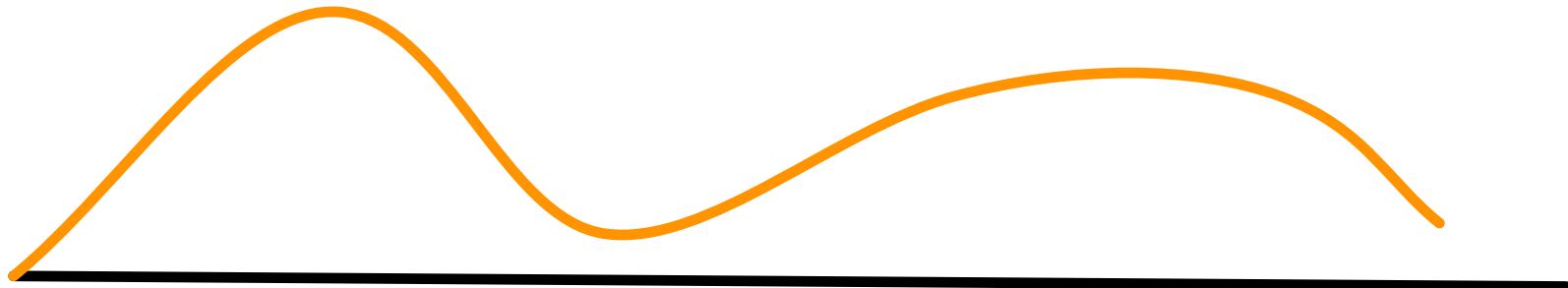
- At test time expected risk matters
(all the other data we could have seen)

$$\mathbf{E}_{(x,y) \sim p} [l(f(x, w), y)]$$



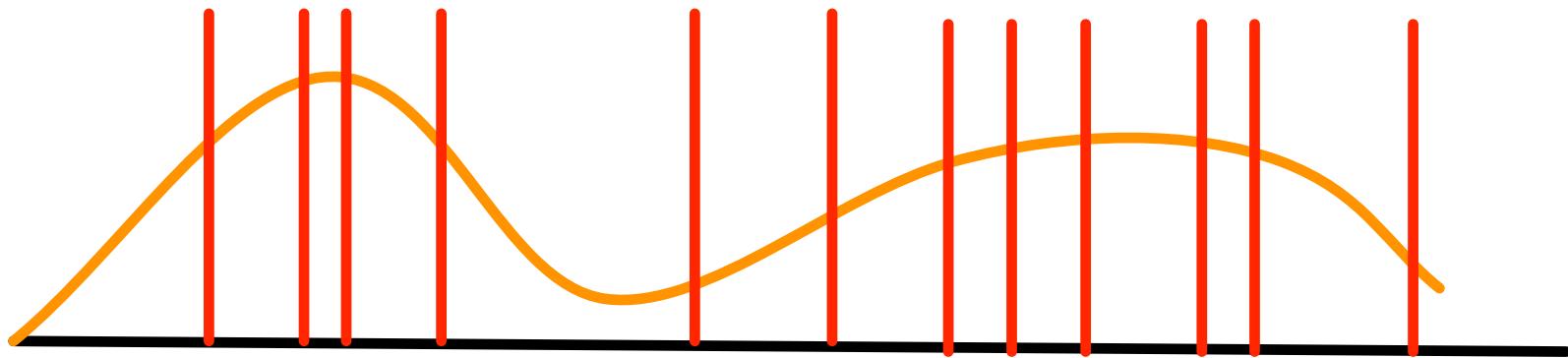
Why

Data Distribution



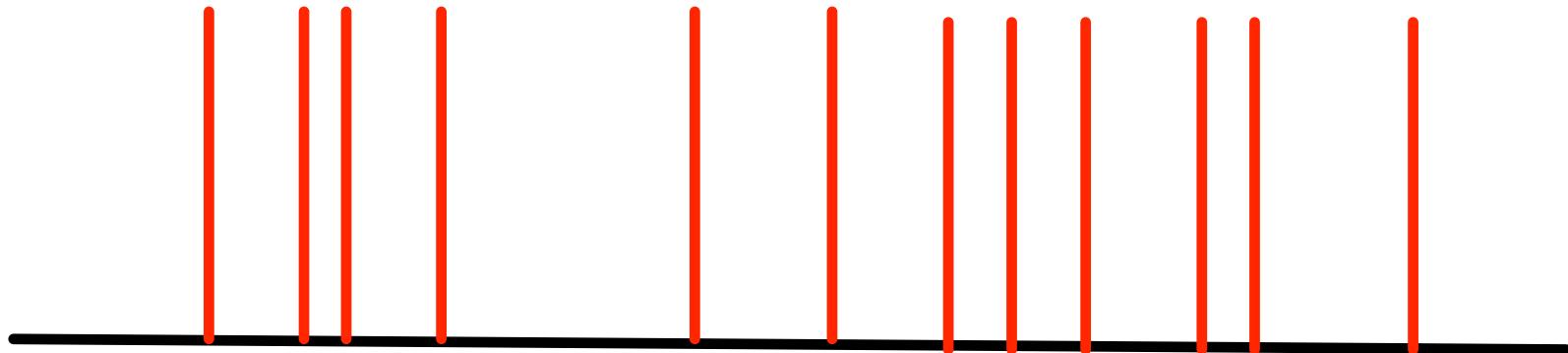
Why

Data Distribution with Empirical Sample



Why

Empirical Sample



Fixing it

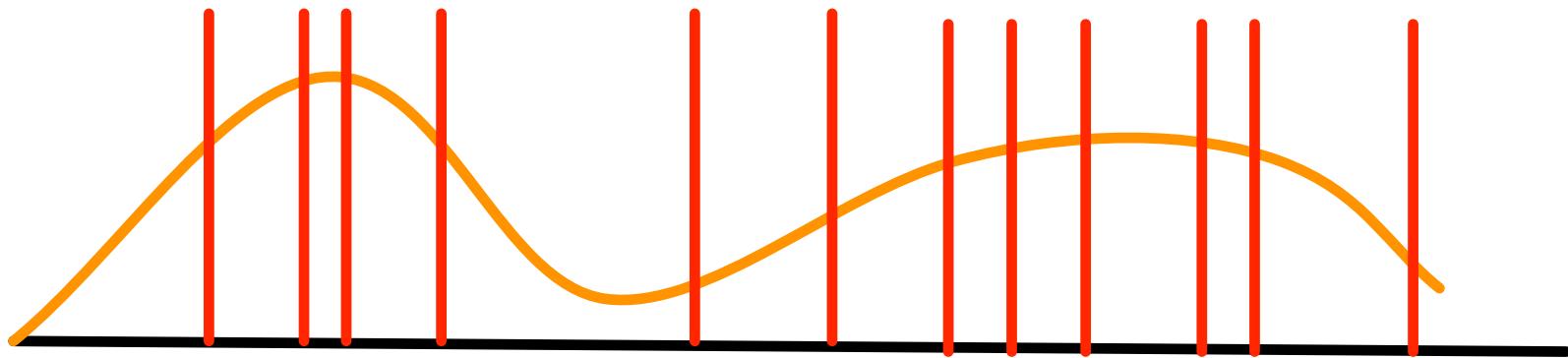
- **Validation set**
(hold out **separate** data that is not used for training)
- **Chernoff bound**

$$\Pr \left\{ \frac{1}{m} \sum_{I=1}^m l(f(x_i), y_i) - \mathbf{E} [l(f(x), y)] > \epsilon \right\} \leq \exp(-2m\epsilon^2)$$

- **Why does it work?**
 - Validation set was never used for training
(often violated)
 - Loss bounded within $[0, 1]$ (otherwise rescale)

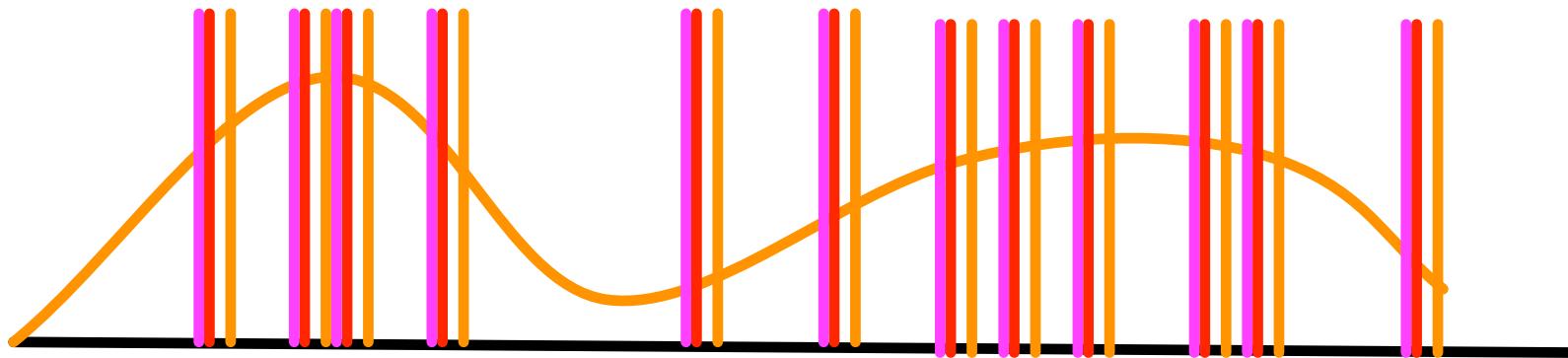
Fixing it

Data Distribution with Empirical Sample



Fixing it

- Input noise (more on this later)
- Dropout (noise within the layers)
- Smoothing the function f (e.g. weight decay)



return

?

/

covariate shift

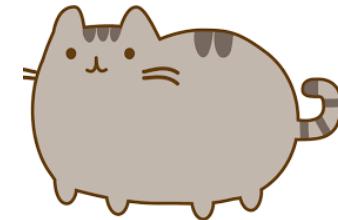
Training set



Training set



At test time



aws

Why would anyone do this?

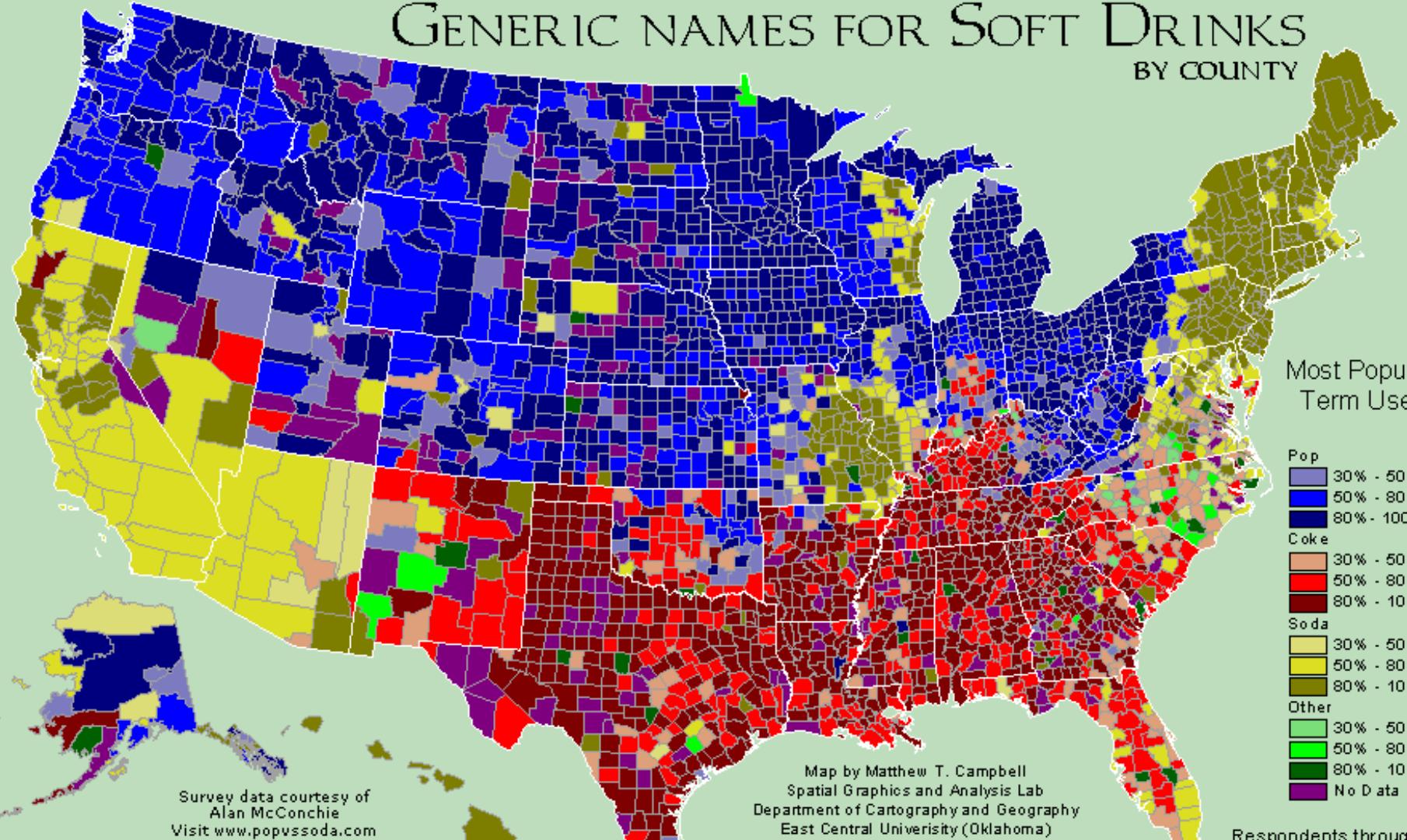


Covariate Shift

- **Web search**
 - Training - page relevance data for the US market
 - Testing - recommend pages for Canada (UK, Australia)
- **Speech recognition**
 - Training - West coast accent
 - Testing - Southern drawl, Texan, non-native speaker
- **Language**
 - Training - ‘James, bring me a **soda**’
 - Testing - ‘John, bring me a ‘**pop**’ (or coke, etc.)



GENERIC NAMES FOR SOFT DRINKS BY COUNTY



Covariate Shift

- **Medical**

- Training - University students + old men with prostate cancer
- Testing - Potentially sick old men

- **Reinforcement Learning**

- Training - Data gathered with current policy
- Testing - Environment reacting to updated policy

- **Databases**

- Training - DB tuned to 2017 usage pattern
- Testing - DB deployed on AWS in 2018



What is happening? $q(x, y) = q(x)p(y|x)$

- Training Risk

Training data

$$\underset{w}{\text{minimize}} \int dx p(x) \int dy p(y|x) l(f(x, w), y)$$

or rather $\underset{w}{\text{minimize}} \frac{1}{m} \sum_{I=1}^m l(f(x_i, w), y_i)$

- Test Risk is different

Test data

$$\int dx q(x) \int dy p(y|x) l(f(x, w), y)$$



Fixing it (covariate shift correction)

- Basic algebra

$$\int dx q(x) f(x) = \int dx p(x) \underbrace{\frac{q(x)}{p(x)}}_{\alpha(x)} f(x) = \int dx p(x) \alpha(x) f(x)$$


- Need to find density ratio, but we don't have either one.
- Estimating p and q directly is really hard and requires specialized tools. Can we recycle classifiers?

Fairness and Bias (covariate shift in action & news)

Training set



cat

dog



dog



dog

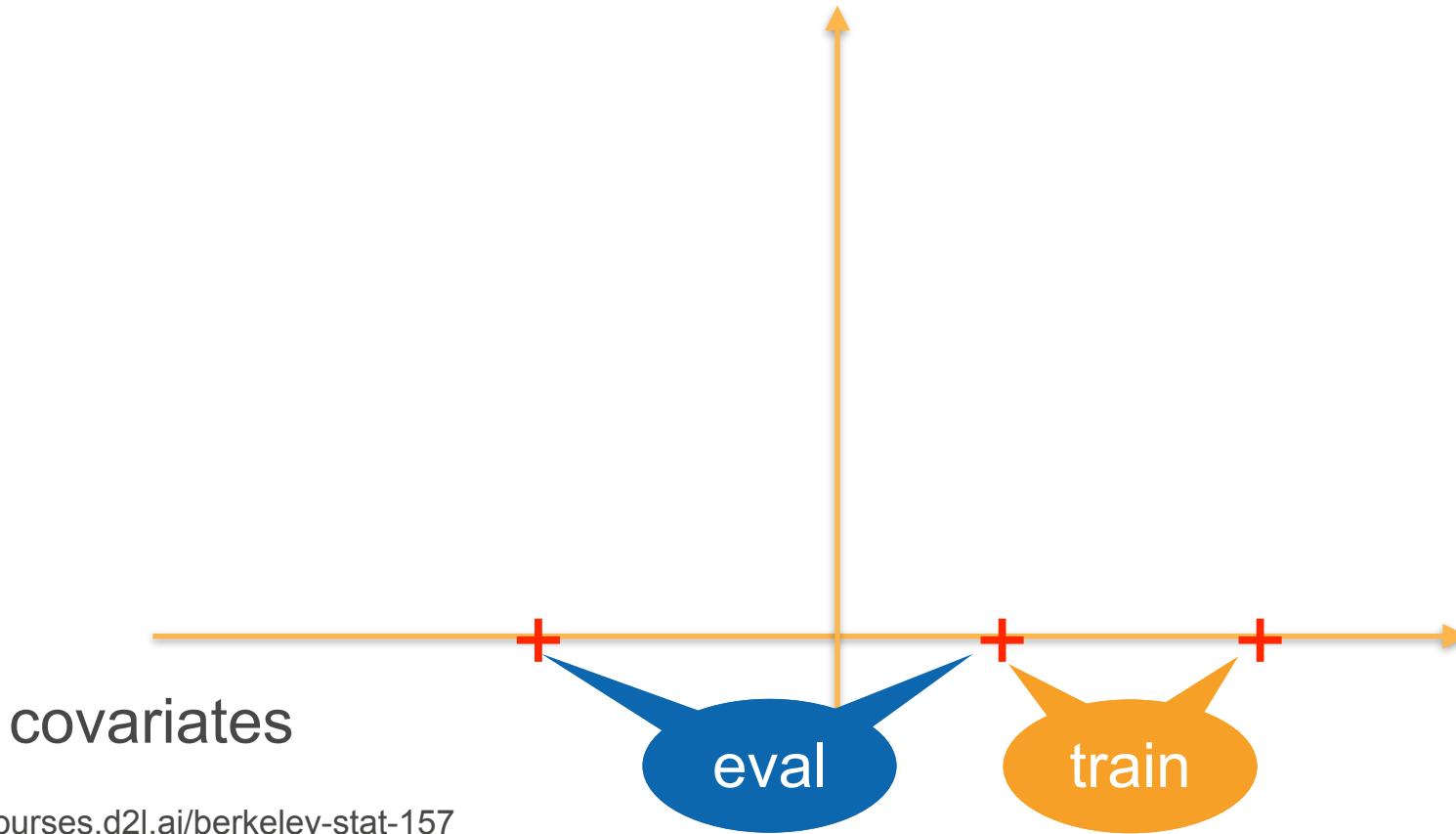
Test set



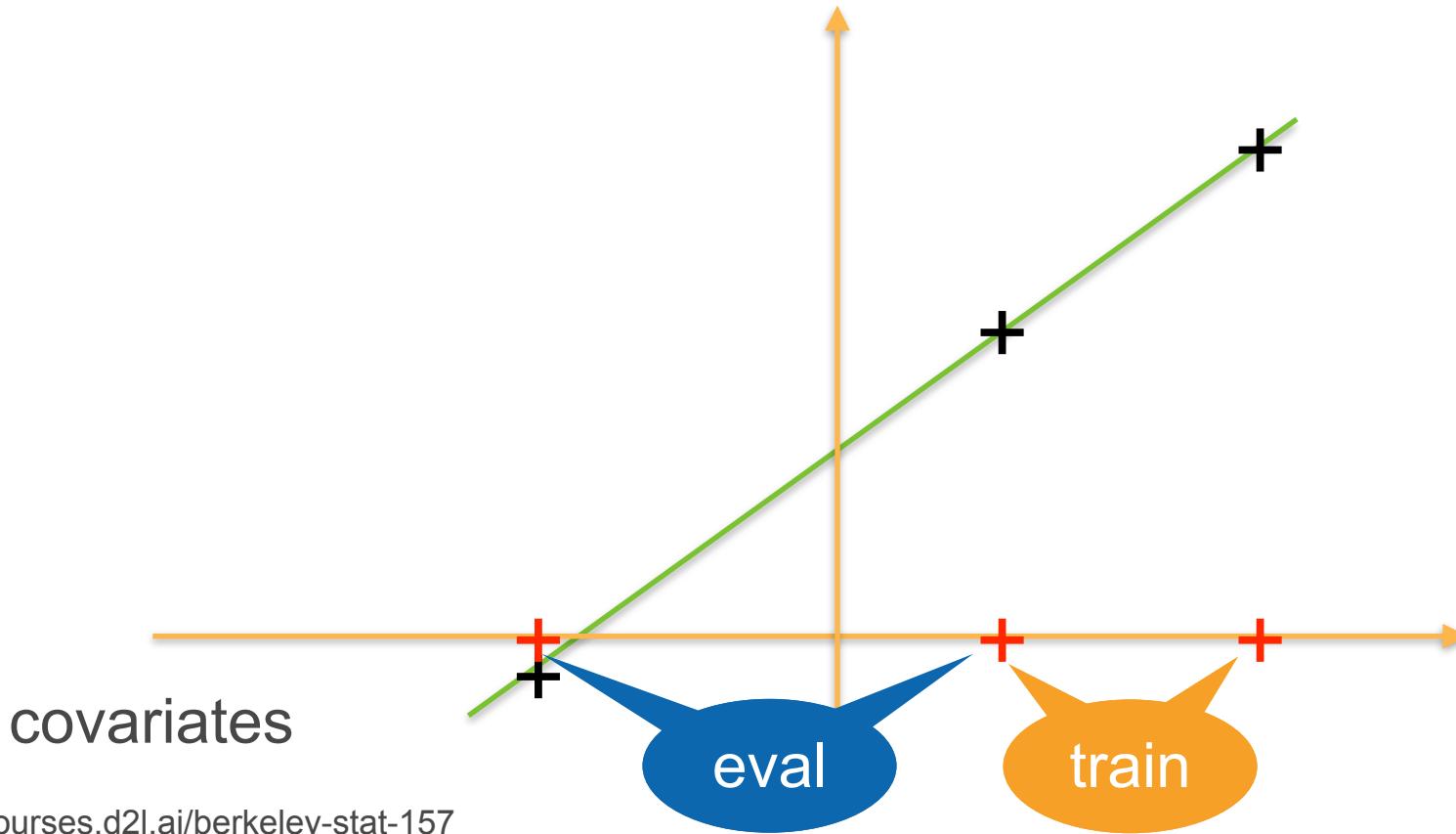
The classifier might perform a lot worse during test time

Why?

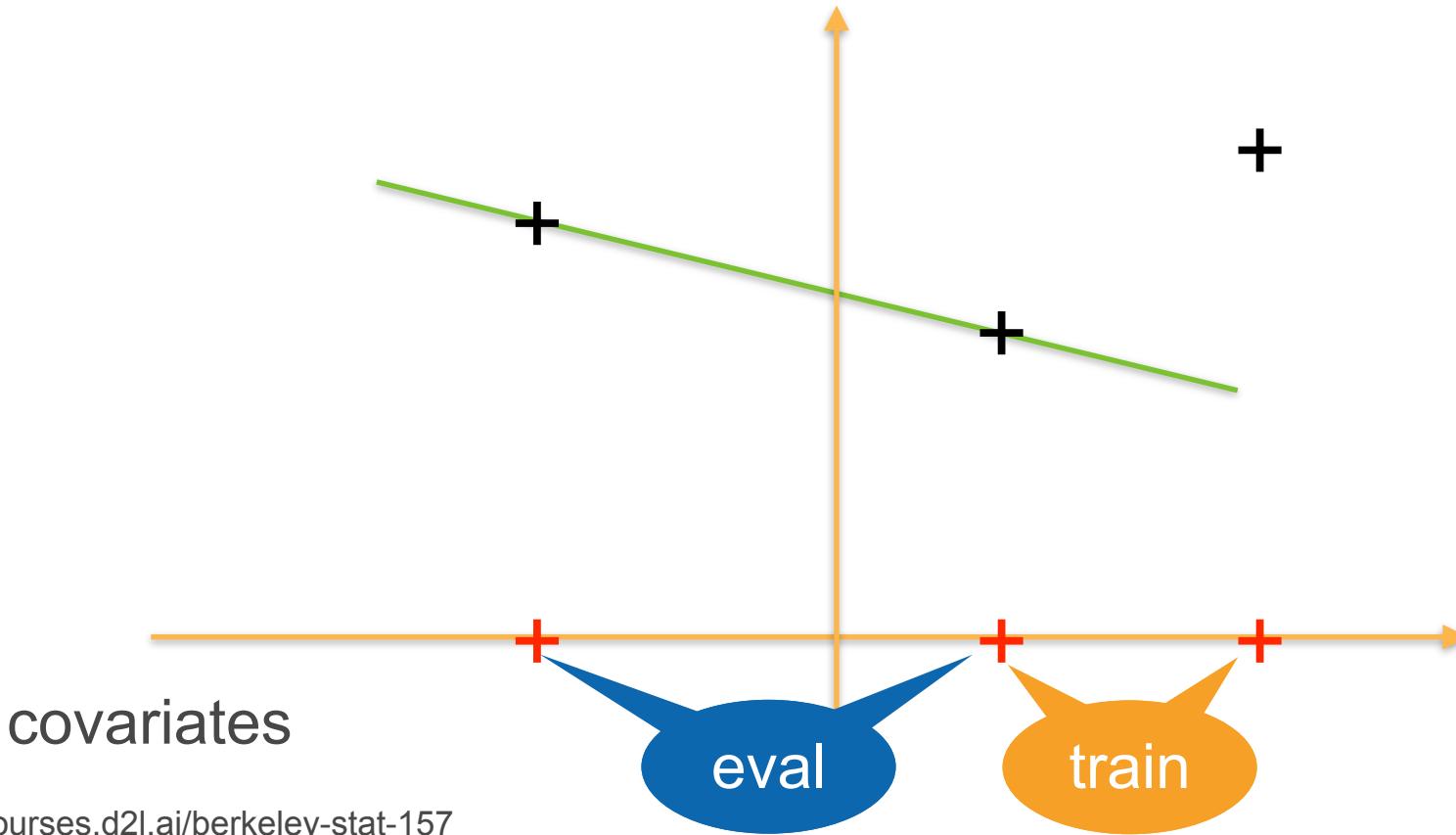
Simple regression problem



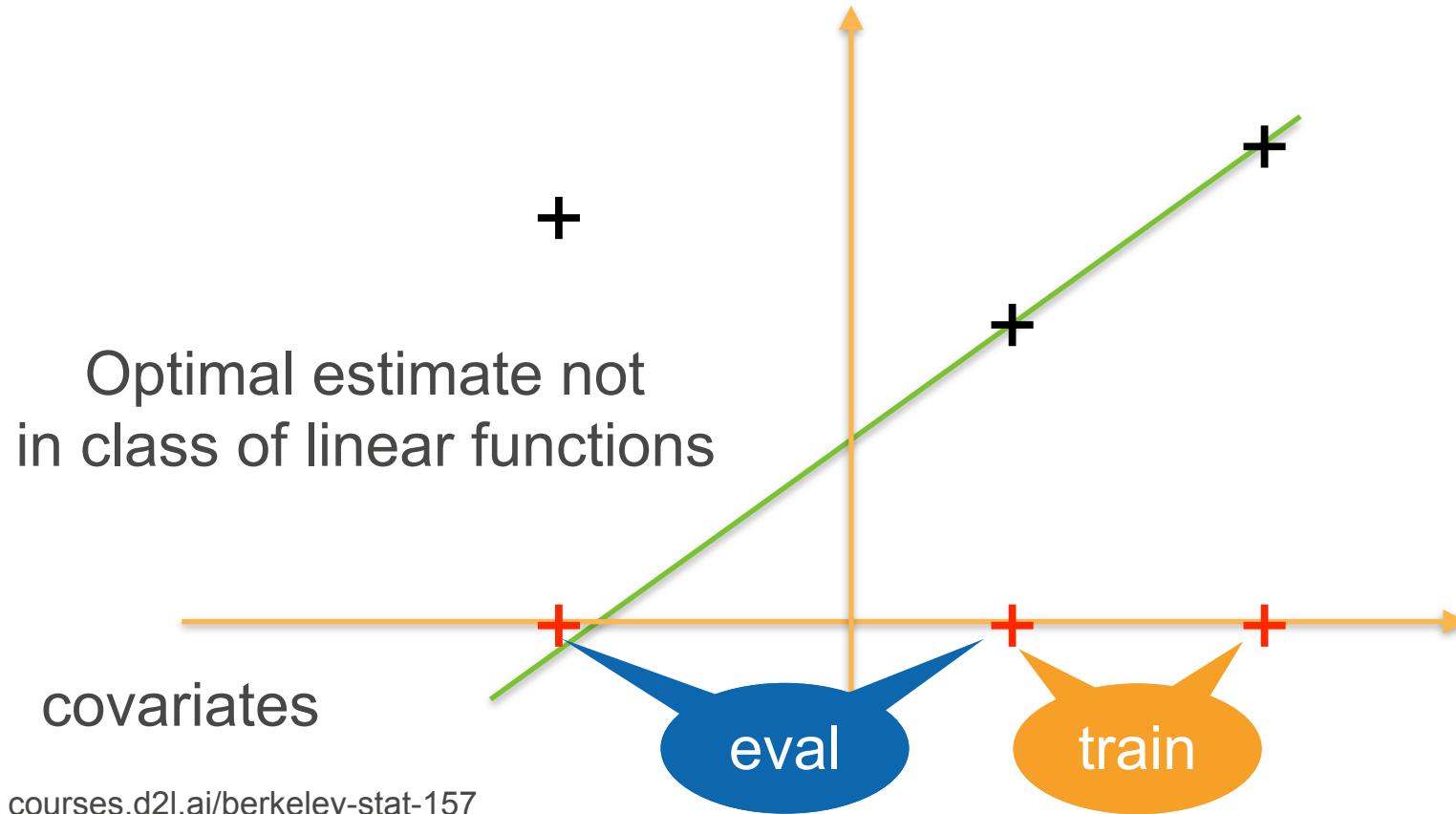
Simple regression problem



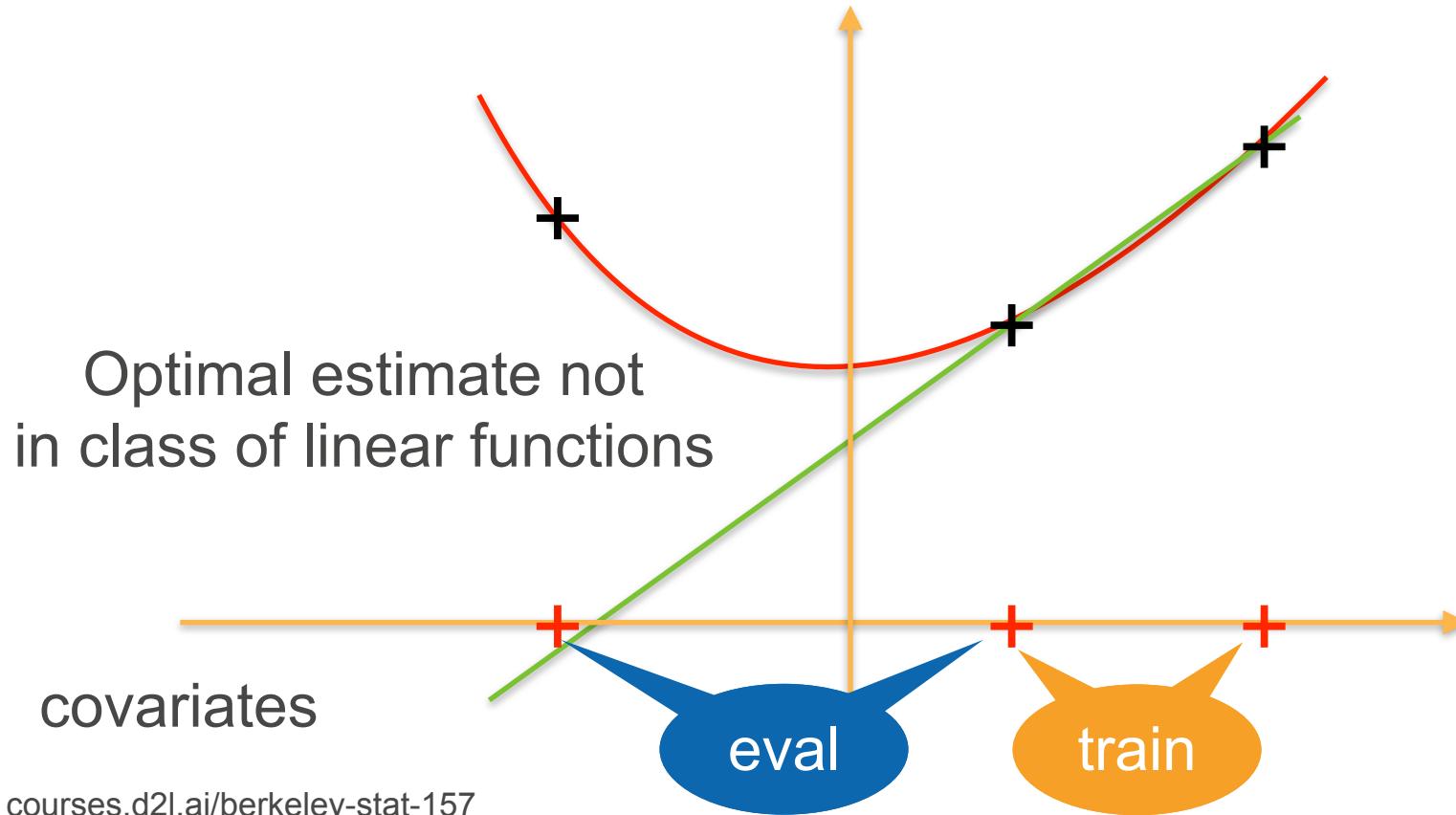
Simple regression problem



Simple regression problem



Simple regression problem



Training error may be misleading (e.g. faces)



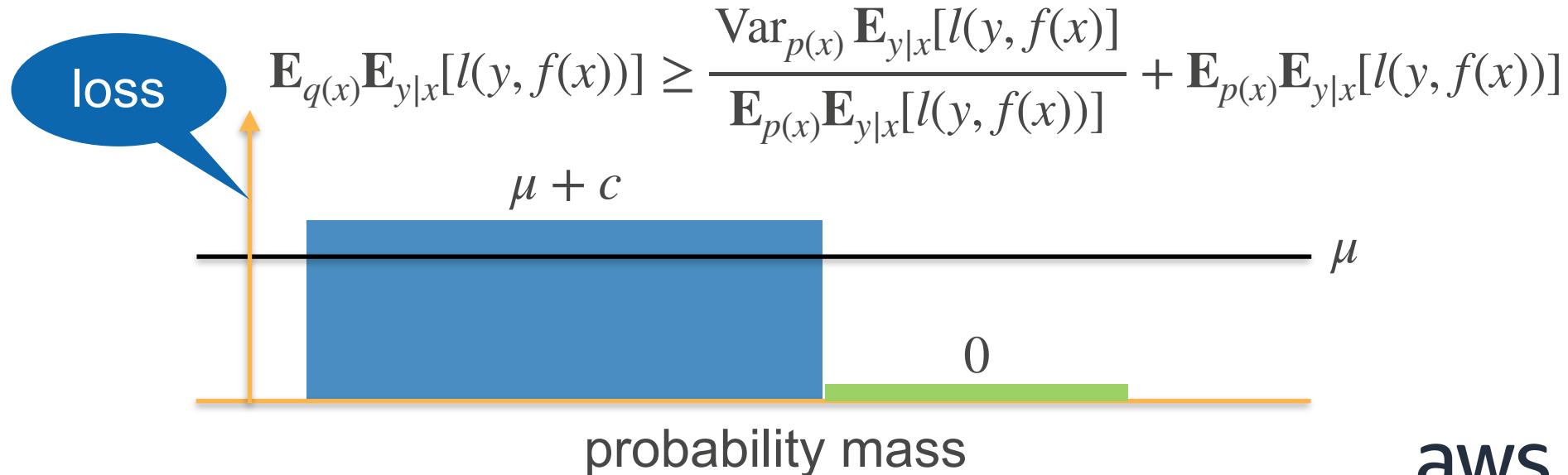
Train on IMDB



Test on weird prof

No Protection against Bias Theorem

- Estimator performs better (or worse) on some data
- We can always find a distribution q that is much worse



TL;DR Testing where we have insufficient amounts of training data *may* yield strange results.



Recall - Multiclass Classification

Calibrated Scale

- Output matches probabilities (nonnegative, sums to 1)

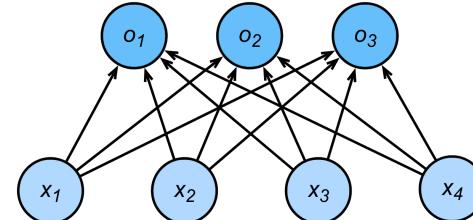
$$\begin{aligned} p(y|o) &= \text{softmax}(o) \\ &= \frac{\exp(o_y)}{\sum_i \exp(o_i)} \end{aligned}$$

- Negative log-likelihood

$$-\log p(y|o) = \log \sum_i \exp(o_i) - o_y$$

Classification

- Multiple classes, typically multiple outputs
- Score *should* reflect confidence ...



Two Classes

- Classes 1 and -1

$$p(y = 1 | o) = \text{softmax}(o) = \frac{\exp(o_1)}{\exp(o_{-1}) + \exp(o_1)}$$

- Shift invariance $o_i \leftarrow o_i + c$

$$p(y = 1 | o) = \frac{\exp(o_1 + c)}{\exp(o_{-1} + c) + \exp(o_1 + c)} = \frac{\exp(o_1)}{\exp(o_{-1}) + \exp(o_1)}$$

Two Classes

- Choose $o_{-1} = 0$

$$p(y = 1 | o) = \frac{\exp(o_1)}{\exp(0) + \exp(o_1)} = \frac{1}{1 + \exp(-o_1)}$$

- Negative log-likelihood

$$-\log p(y | o) = \log(1 + \exp(-yo_1))$$

The graph illustrates the behavior of the logistic loss function as x approaches positive and negative infinity. Three curves are shown: a green curve that decreases from approximately 5 at $x = -5$ towards 0; a blue curve that decreases from approximately 4 at $x = -5$ towards 0; and an orange curve that increases from approximately -1 at $x = -5$ towards 0. All three curves approach 0 as $|x|$ increases.

$$\lim_{x \rightarrow \infty} \log(1 + \exp(-x)) = \log 1 = 0$$
$$\lim_{x \rightarrow -\infty} \log(1 + \exp(-x)) + x = \lim_{x \rightarrow -\infty} \log(1 + \exp(x)) = 0$$

Logistic loss function



Logistic Regression Summary

- **Data** (x_i, y_i) where $x_i \in \mathcal{X}$ and $y_i \in \{\pm 1\}$
- **Objective**
$$\underset{w}{\text{minimize}} - \sum_{i=1}^m \log(1 + \exp(-y_i f(x_i, w))) + \text{penalty}(w)$$

- **Conditional Probability Estimate**

$$\log p(y = 1 | o) = \frac{1}{1 + \exp(-o)}$$

return

covariate shift correction

Covariate shift correction

- Propensity scoring

$$\int dx q(x) f(x) = \int dx p(x) \underbrace{\frac{q(x)}{p(x)}}_{\alpha(x)} f(x) = \int dx p(x) \alpha(x) f(x)$$


- Need to find density ratio, but we don't have either one.
- Key idea: train a classifier between p and q

$$r(x, y) = \frac{1}{2} [p(x)\delta(y, 1) + q(x)\delta(y, -1)]$$



Covariate shift correction

- Conditional class probability

$$r(y = 1|x) = \frac{p(x)}{p(x) + q(x)} \text{ and hence } \alpha = \frac{q(x)}{p(x)} = \frac{r(y = -1|x)}{r(y = 1|x)}$$

- Logistic regression

$$r(y = 1 | x) = \frac{1}{1 + \exp(-f(x))}$$

$$\implies \alpha(x) = \frac{r(y = -1 | x)}{r(y = 1 | x)} = \exp(f(x))$$



Covariate Shift Correction Redux

- Training and test data
- Split **as if it were a binary classification problem** (labels -1 and 1 for training and test respectively)
- Train with **logistic regression** to get f
- Use binary classifier output to reweight data
- Solve original problem but weighted

$$\sum_i l(x_i, y_i, g(x_i, w)) \longrightarrow \sum_i \exp(f(x_i)) \cdot l(x_i, y_i, g(x_i, w))$$

Label shift

Training set



Training set



Test set



A vertical orange line separates the cat images from the dog images.



Why would anyone do this?



Label Shift

- **Medical diagnosis**
 - Train on data with few sick patients
 - Test on data during flu season where $q(\text{flu}) > p(\text{flu})$ while flu symptoms $p(\text{symptoms}|\text{flu})$ are still the same
- **Speech recognition**
 - Train on newscast data before election
 - Test on newscast after election (new topics, names, discussions, but still same language)

Label Shift

$$q(x, y) = q(y)p(x|y)$$

- Data generating process $p(x|y)$ is unchanged
- Labels change since the underlying cause changed
- Need to reweight according to $\beta(y) = \frac{q(y)}{p(y)}$ to get

$$\int q(y) dy \int p(x|y) dx l(f(x), y) = \int p(y) \frac{q(y)}{p(y)} dy \int p(x|y) dx l(f(x), y)$$



We don't have samples from $q(y)!$



Label Shift

$$q(x, y) = q(y)p(x|y)$$

- **Key Idea - measure the estimates on test set**
 - $p(x|y)$ is the same for training and test
 - Distribution of predictions $|x, y$ has to be the same
- **Simple ‘spectral’ algorithm** (Lipton, Wang, Smola, 2018)
 - Confusion matrix $C[y'|y] = \Pr(\hat{y}(x) = y'|y)$ on hold out
 - Predicted label vector on test set $\mu[y'] = \Pr(\hat{y}(x) = y')$
 - Obtain $q(y)$ via matrix inversion since

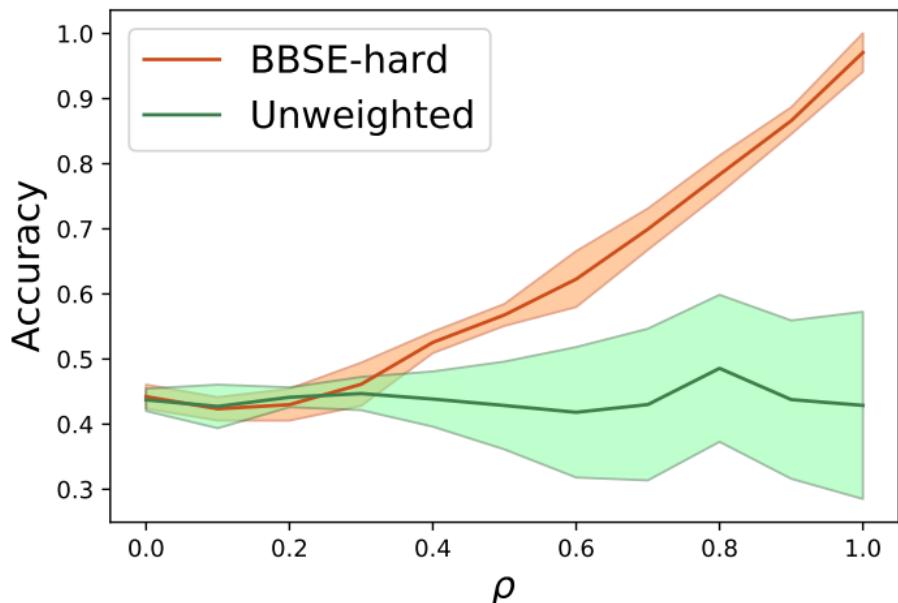
$$\mu[y'] = \sum_y C[y'|y]q(y)$$



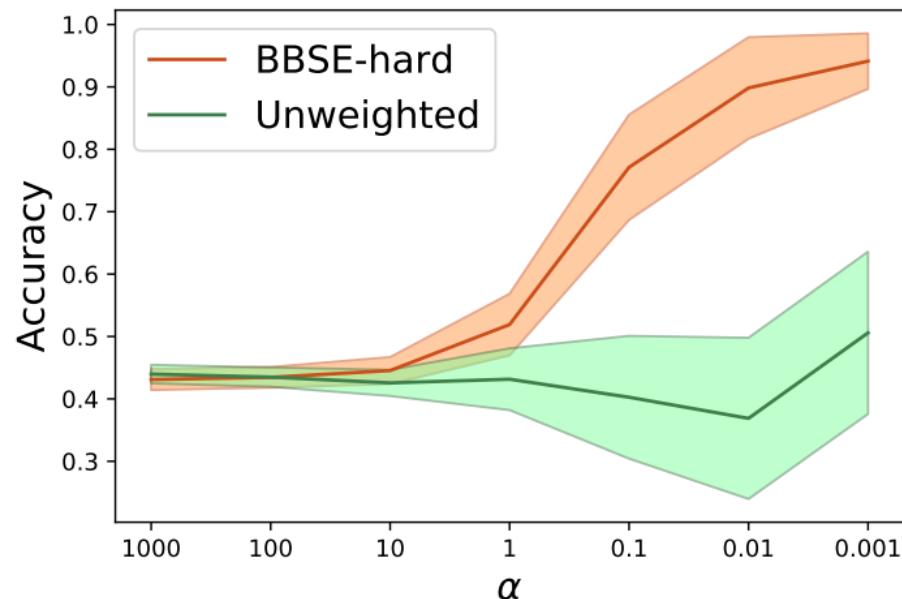
Guarantees

- **Robust under misspecification**
 - Even if the estimates $y(x)$ are wrong, calibration is OK:
(same errors on hold-out and test set)
 - Confusion matrix and label vector are concentrated:
(use matrix Bernstein inequality)
- **Simple algorithm**
 - Cubic in number of classes, linear in sample size

Black Box Shift Correction on CIFAR10



Tweaking one class probability

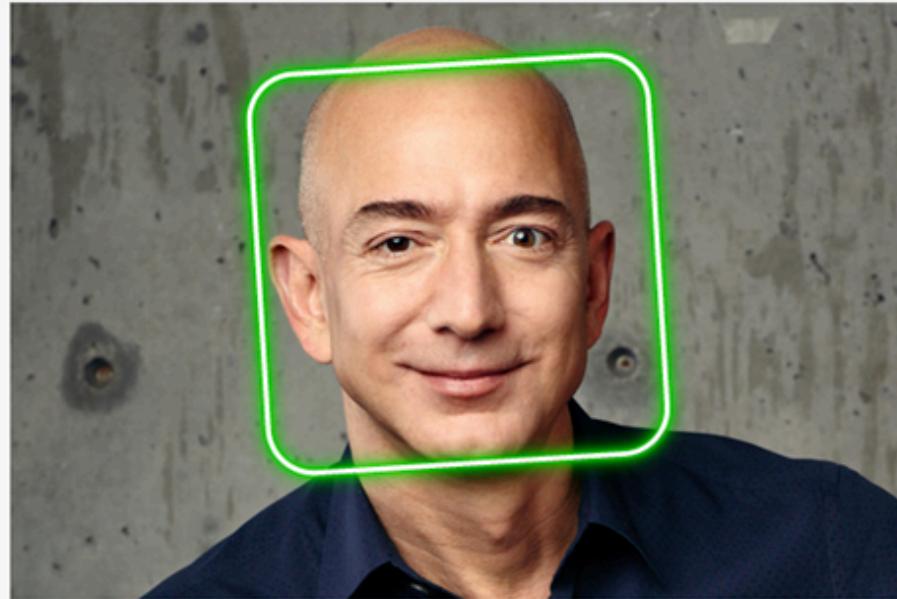


Dirichlet prior over shifts

Adversarial data

Celebrity recognition

Rekognition automatically recognizes celebrities in images and provides confidence scores (Your images aren't stored.)



Choose a sample Image



Use your own image

 Upload

or drag and drop

Use image URL

Go

Done with the demo?

[Download SDKs](#)

▼ Results



Jeff Bezos
[Learn More](#)

Match confidence

100%

► Request

► Response

aws

Adversarial Image Generation (e.g. Sharif et al. 2017)

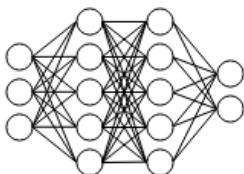
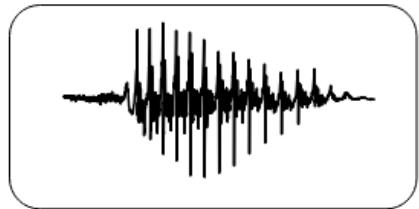


Digital manipulation
to dodge recognition



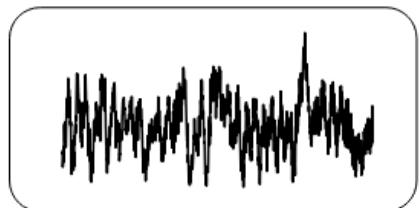
In real life - via 3D
printed glasses

Adversarial Audio Generation (e.g. Carlini & Wagner, 2018)



"it was the
best of times,
it was the
worst of times"

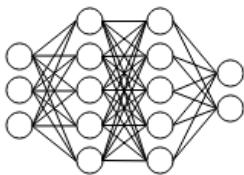
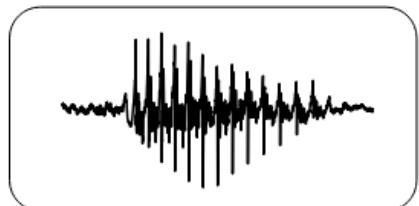
- Modify data slightly such as to obtain wrong class



+

$$\times \ 0.001$$

=



"it is a truth
universally
acknowledged
that a single"

$$\underset{\delta}{\text{maximize}} \ l(f(x + \delta), y)$$

$$\text{subject to } \|\delta\| \leq \epsilon$$

Different norms
Different datasets
Different papers ...

Why does this work?



'Unnatural' data



- Training and 'natural' test data live in small subset
- Adversarial data is slightly off that support
- Function behavior undefined away from where data occurs

'Unnatural' data



- Training and 'natural' test data live in small subset
- Adversarial data is slightly off that support
- Function behavior undefined away from where data occurs

Wow. Breathtaking. Is this new?

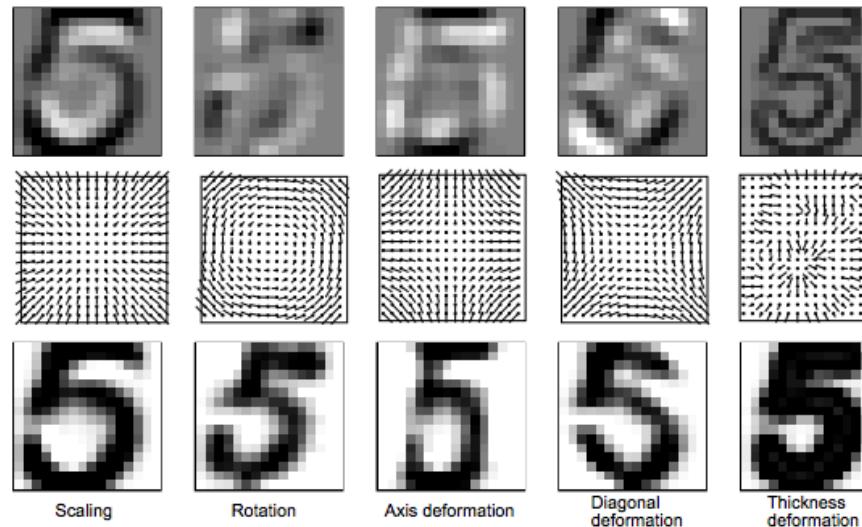


Spam defenses

- **While TRUE**
 - Mail host extends dataset and trains new classifier
 - Spammer's e-mails are rejected
 - Spammer finds a modification that succeeds
- **Examples**
 - Add highly scoring words (or sentences) to email
 - Add highly scoring sentences (and vary them)
 - Change or forge header ('Dear Alex, ...')

Invariances

- Tangent Distance (Simard et al., 1995)
 - Invariance transforms don't change the label
 - Explore data and their neighborhood



Invariances

- **Virtual Support Vectors** (Schoelkopf, 1997)
Only change the data at the boundary (not enough RAM)
- **Data augmentation for training**
 - **Imagenet** (pretty much every paper)
Cropping, scaling, change mean, per channel, ...
 - **Speech Recognition**
Background noise, scenes, ...
 - **Document Analysis**
Random substrings, word removal, insertion



Invariant and robust loss

- **Convex loss** (Teo et al, 2005)
 - Family of transformations $\delta \in \Delta$
 - Penalty for extreme transformations $1 \geq \eta(\delta) \geq 0$
 - Find the ‘worst’ possible example at each step

Adversarially Robust
Networks

$$L(x, y, f) = \sup_{\delta \in \Delta} \eta(\delta) l(f(x + \delta), y)$$

e.g. adversarial
example generator
Finds worst possible

Reduced penalty for
extreme distortions

Nonstationary Environments

Interaction with Environment

- **Batch** (download a book)
Observe training data $(x_1, y_1) \dots (x_l, y_l)$ then deploy
- **Online** (follow the class)
Observe x , predict $f(x)$, observe y (stock market, homework)
- **Active learning** (ask questions in class)
Query y for x , improve model, pick new x
- **Bandits** (do well at homework)
Pick arm, get reward, pick new arm (also with context)
- **Reinforcement Learning** (play chess, drive a car)
Take action, environment responds, take new action

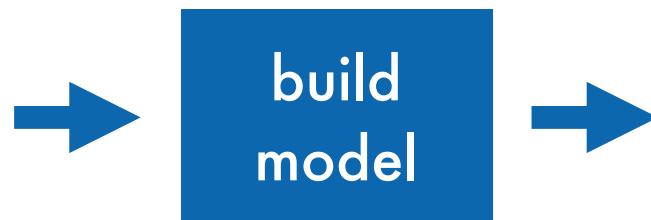
Batch

training data

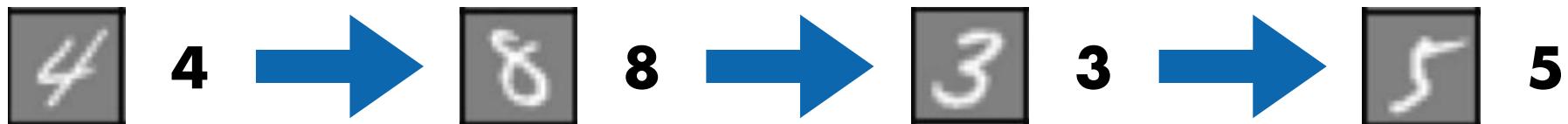
6	5	5	4	1	0
7	4	0	8	4	3
3	4	2	8	1	0
0	0	1	6	5	5
1	1	1	6	7	1
8	6	4	5	3	8
1	7	2	8	4	7
5	2	8	0	4	8
3	3	7	0	5	3
4	8	9	4	0	4

test data

4	9	1	7
6	4	5	6
7	5	9	7
1	1	5	9
4	1	3	1
7	2	9	1
6	8	9	3
3	7	+	6
1	1	0	3
5	0	5	0



Online



System improves as we see more data

Bandits

- Choose an arm (action)
- See what happens (get reward)
- Update model
- Choose next arm (action)

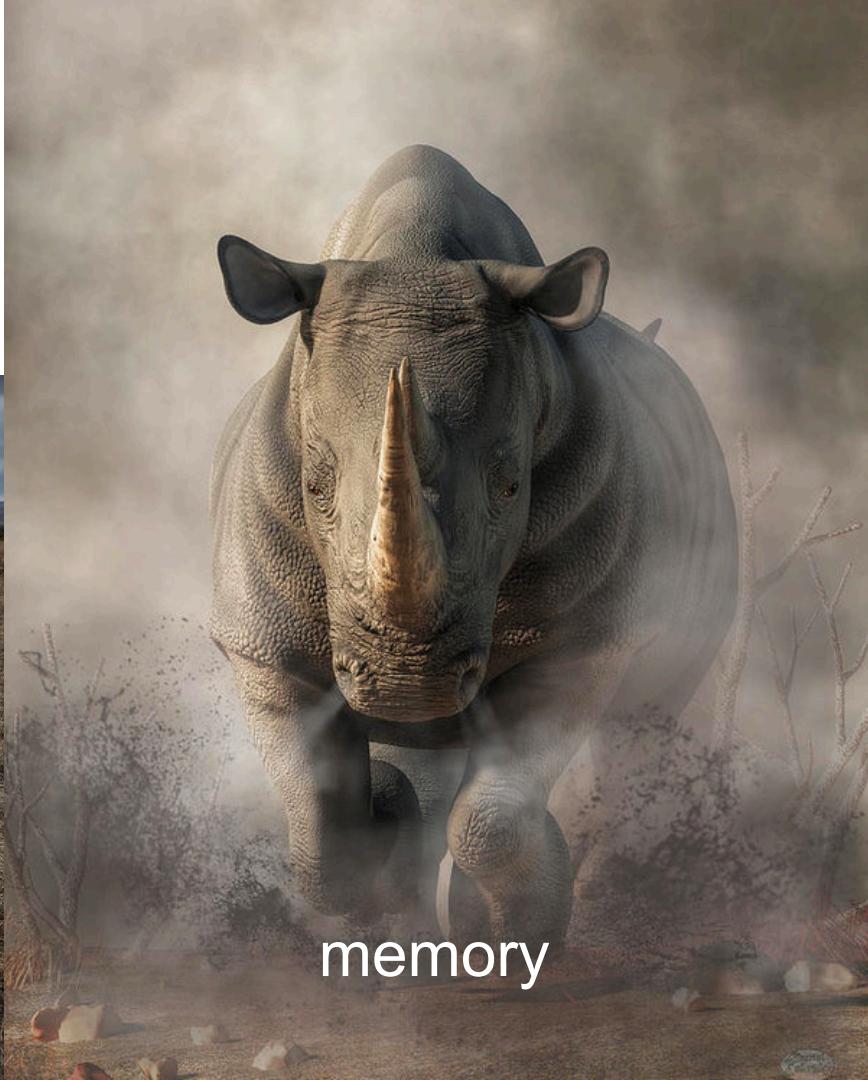
The bandit **doesn't remember** what you did last summer.



Stateful Systems



no memory



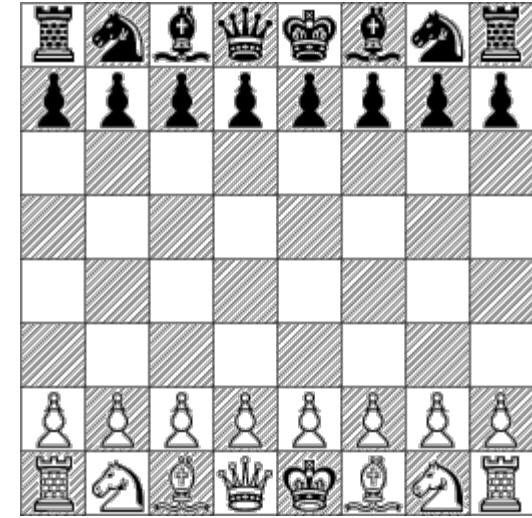
memory

Reinforcement Learning & Control

- Take action
- Environment reacts
- Observe stuff
- Update model

Repeat

- environment (cooperative, adversary, doesn't care)
- memory (goldfish, elephant)
- state space (tic tac toe, chess, car)
- past observations (server log, generated during training)



Reinforcement Learning & Control

- **Games**

- Chess, Go, Backgammon (fully observed)
- Poker, Starcraft, ATARI (partially observed, random)

- **Parallelism**

- Computation advertising, recommender systems (multiple agents & independent parallel games)
- Load balancing & scheduling (multiple agents)

- **Actions**

- Continuous decisions (driving, flying, robots in general, HVAC)
- Discrete (elevator, work allocation)

- **Simulations**

- MuJoCo style
- Only reality (server center)

Training ≠ Testing

- **Generalization performance**
(the empirical distribution lies)
- **Covariate shift**
(the covariate distribution lies)
- **Logistic regression**
(tools to fix shift)
- **Covariate shift correction**
- **Label shift**
(the label distribution lies)
- **Nonstationary Environments**

$$p_{\text{emp}}(x, y) \neq p(x, y)$$

$$p(x) \neq q(x)$$

$$\log(1 + \exp(-yf(x)))$$

$$\frac{1}{2} (p(x)\delta(1, y) + q(x)\delta(-1, y))$$

$$p(y) \neq q(y)$$



Numerical Stability



Wikipedia

Gradients for Neural Networks

- Consider a network with d layers

$$\mathbf{h}^t = f_t(\mathbf{h}^{t-1}) \quad \text{and} \quad y = \ell \circ f_d \circ \dots \circ f_1(\mathbf{x})$$

- Compute the gradient of the loss ℓ w.r.t. \mathbf{W}_t

$$\frac{\partial \ell}{\partial \mathbf{W}^t} = \frac{\partial \ell}{\partial \mathbf{h}^d} \underbrace{\frac{\partial \mathbf{h}^d}{\partial \mathbf{h}^{d-1}} \cdots \frac{\partial \mathbf{h}^{t+1}}{\partial \mathbf{h}^t}}_{\text{Multiplication of } d-t \text{ matrices}} \frac{\partial \mathbf{h}^t}{\partial \mathbf{W}^t}$$

Multiplication of $d-t$ matrices

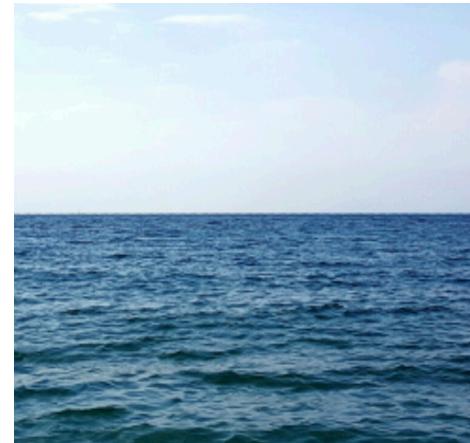
Two Issues for Deep Neural Networks

$$\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i}$$

Gradient Exploding



Gradient Vanishing



$$1.5^{100} \approx 4 \times 10^{17}$$

$$0.8^{100} \approx 2 \times 10^{-10}$$

Example: MLP

- Assume MLP (without bias for simplicity)

$$f_t(\mathbf{h}^{t-1}) = \sigma(\mathbf{W}^t \mathbf{h}^{t-1}) \quad \sigma \text{ is the activation function}$$

$$\frac{\partial \mathbf{h}^t}{\partial \mathbf{h}^{t-1}} = \text{diag} (\sigma'(\mathbf{W}^t \mathbf{h}^{t-1})) (\mathbf{W}^t)^T \quad \sigma' \text{ is the gradient function of } \sigma$$

$$\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i} = \prod_{i=t}^{d-1} \text{diag} (\sigma'(\mathbf{W}^i \mathbf{h}^{i-1})) (\mathbf{W}^i)^T$$

Gradient Exploding

- Use ReLU as the activation function

$$\sigma(x) = \max(0, x) \quad \text{and} \quad \sigma'(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{otherwise} \end{cases}$$

- Elements of $\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i} = \prod_{i=t}^{d-1} \text{diag}(\sigma'(\mathbf{W}^i \mathbf{h}^{i-1})) (\mathbf{W}^i)^T$ may from $\prod_{i=t}^{d-1} (\mathbf{W}^i)^T$
 - Leads to large values when $d-t$ is large

$$1.5^{100} \approx 4 \times 10^{17}$$

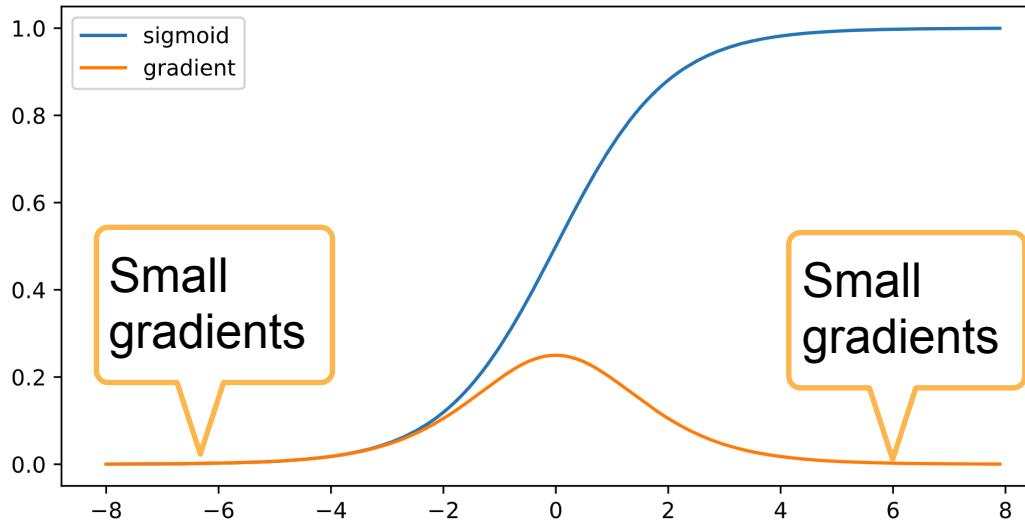
Issues with Gradient Exploding

- Value out of range: infinity value
 - Severe for using 16-bit floating points
 - Range: $6e-5 - 6e4$
- Sensitive to learning rate (LR)
 - Not small enough LR \rightarrow large weights \rightarrow larger gradients
 - Too small LR \rightarrow No progress
 - May need to change LR dramatically during training

Gradient Vanishing

- Use sigmoid as the activation function

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \sigma'(x) = \sigma(x)(1 - \sigma(x))$$



Gradient Exploding

- Use sigmoid as the activation function

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad \sigma'(x) = \sigma(x)(1 - \sigma(x))$$

- Elements $\prod_{i=t}^{d-1} \frac{\partial \mathbf{h}^{i+1}}{\partial \mathbf{h}^i} = \prod_{i=t}^{d-1} \text{diag}(\sigma'(\mathbf{W}^i \mathbf{h}^{i-1})) (\mathbf{W}^i)^T$ are products of $d-t$ small values

$$0.8^{100} \approx 2 \times 10^{-10}$$

Issues with Gradient Vanishing

- Gradients with value 0
 - Severe with 16-bit floating points
- No progress in training
 - No matter how to choose learning rate
- Severe with bottom layers
 - Only top layers are well trained
 - No benefit to make networks deeper

Stabilize Training

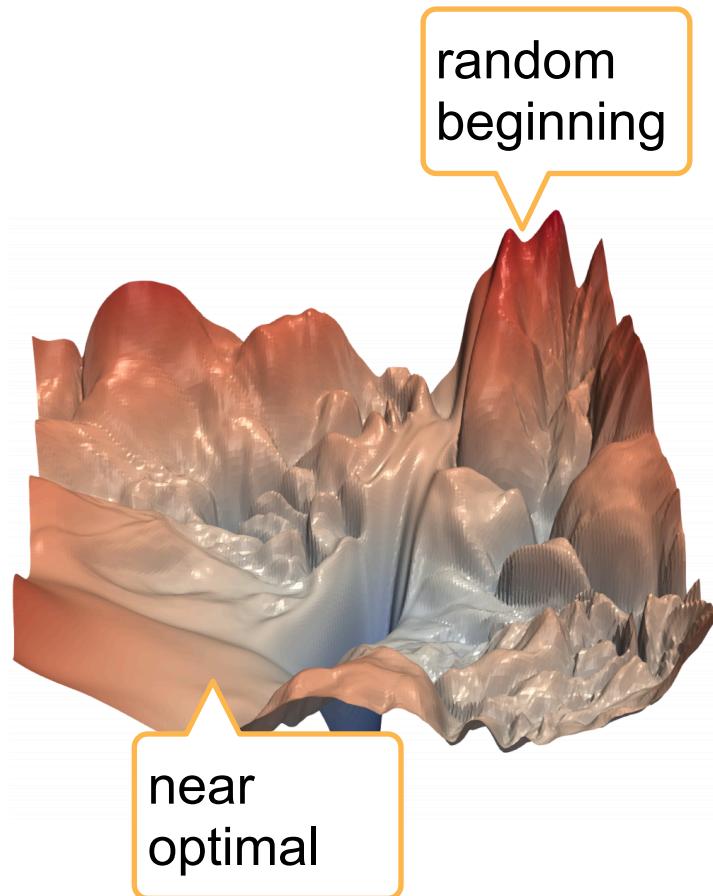


Stabilize Training

- Goal: make sure gradient values are in a proper range
 - E.g. in $[1e-6, 1e3]$
- Multiplication -> plus
 - ResNet, LSTM
- Normalize
 - Batch Normalization, Gradient clipping
- Proper weight initialization and activation functions

Weight Initialization

- Initialize weights with random values in a proper range
- The beginning of training easily suffers to numerical instability
 - The surface far away from an optimal can be complex
 - Near optimal may be flatter
- Initializing according to $\mathcal{N}(0, 0.01)$ works well for small networks, but not guarantee for deep neural networks



Constant Variance for each Layer

- Treat both layer outputs and gradients are random variables
- Make the mean and variance for each layer's output are same, similar for gradients

Forward

$$\begin{aligned}\mathbb{E}[h_i^t] &= 0 \\ \text{Var}[h_i^t] &= a\end{aligned}\quad \mathbb{E} \left[\frac{\partial \ell}{\partial h_i^t} \right] = 0 \quad \text{Var} \left[\frac{\partial \ell}{\partial h_i^t} \right] = b \quad \forall i, t$$

a and b are constants

Example: MLP

- Assumptions

- i.i.d $w_{i,j}^t$, $\mathbb{E}[w_{i,j}^t] = 0$, $\text{Var}[w_{i,j}^t] = \gamma_t$
- h_i^{t-1} is independent to $w_{i,j}^t$
- identity activation: $\mathbf{h}^t = \mathbf{W}^t \mathbf{h}^{t-1}$ with $\mathbf{W}^t \in \mathbb{R}^{n_t \times n_{t-1}}$

$$\mathbb{E}[h_i^t] = \mathbb{E} \left[\sum_j w_{i,j}^t h_j^{t-1} \right] = \sum_j \mathbb{E}[w_{i,j}^t] \mathbb{E}[h_j^{t-1}] = 0$$

Forward Variance

$$\begin{aligned}\text{Var}[h_i^t] &= \mathbb{E}[(h_i^t)^2] - \mathbb{E}[h_i^t]^2 = \mathbb{E} \left[\left(\sum_j w_{i,j}^t h_j^{t-1} \right)^2 \right] \\ &= \mathbb{E} \left[\sum_j \left(w_{i,j}^t \right)^2 \left(h_j^{t-1} \right)^2 + \sum_{j \neq k} w_{i,j}^t w_{i,k}^t h_j^{t-1} h_k^{t-1} \right] \quad \Rightarrow \quad n_{t-1} \gamma_t = 1 \\ &= \sum_j \mathbb{E} \left[\left(w_{i,j}^t \right)^2 \right] \mathbb{E} \left[\left(h_j^{t-1} \right)^2 \right] \\ &= \sum_j \text{Var}[w_{i,j}^t] \text{Var}[h_j^{t-1}] = n_{t-1} \gamma_t \text{Var}[h_j^{t-1}]\end{aligned}$$

Backward Mean and Variance

- Apply forward analysis as well

$$\frac{\partial \ell}{\partial \mathbf{h}^{t-1}} = \frac{\partial \ell}{\partial \mathbf{h}^t} \mathbf{W}^t \text{ leads to } \left(\frac{\partial \ell}{\partial \mathbf{h}^{t-1}} \right)^T = (\mathbf{W}^t)^T \left(\frac{\partial \ell}{\partial \mathbf{h}^t} \right)^T$$

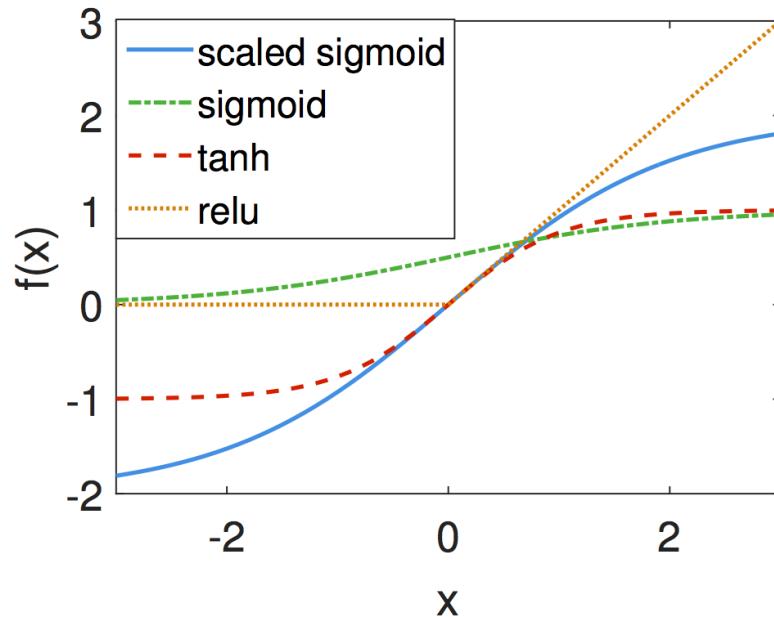
$$\mathbb{E} \left[\frac{\partial \ell}{\partial h_i^{t-1}} \right] = 0$$

$$\text{Var} \left[\frac{\partial \ell}{\partial h_i^{t-1}} \right] = n_t \gamma_t \text{Var} \left[\frac{\partial \ell}{\partial h_j^t} \right] \quad \Rightarrow \quad n_t \gamma_t = 1$$

Xavier Initialization

- Conflict goal to satisfies both $n_{t-1}\gamma_t = 1$ and $n_t\gamma_t = 1$
- Xavier $\gamma_t(n_{t-1} + n_t)/2 = 1 \rightarrow \gamma_t = 2/(n_{t-1} + n_t)$
 - Normal distribution $\mathcal{N}\left(0, \sqrt{2/(n_{t-1} + n_t)}\right)$
 - Uniform distribution $\mathcal{U}\left(-\sqrt{6/(n_{t-1} + n_t)}, \sqrt{6/(n_{t-1} + n_t)}\right)$
 - Variance of $\mathcal{U}[-a, a]$ is $a^2/3$
- Adaptive to weight shape, especially when n_t varies

Activation



A Simple Linear Activation Function

- Assume $\sigma(x) = \alpha x + \beta$

$$\mathbf{h}' = \mathbf{W}^t \mathbf{h}^{t-1} \quad \text{and} \quad \mathbf{h}^t = \sigma(\mathbf{h}')$$

$$\mathbb{E}[h_i^t] = \mathbb{E} [\alpha h'_i + \beta] = \beta \qquad \Rightarrow \qquad \beta = 0$$

$$\begin{aligned}\text{Var}[h_i^t] &= \mathbb{E}[(h_i^t)^2] - \mathbb{E}[h_i^t]^2 \\ &= \mathbb{E}[(\alpha h'_i + \beta)^2] - \beta^2 \qquad \Rightarrow \qquad \alpha = 1 \\ &= \mathbb{E}[\alpha^2(h'_i)^2 + 2\alpha\beta h'_i + \beta^2] - \beta^2 \\ &= \alpha^2 \text{Var}[h'_i]\end{aligned}$$

Backward

- Assume $\sigma(x) = \alpha x + \beta$

$$\frac{\partial \ell}{\partial \mathbf{h}'} = \frac{\partial \ell}{\partial \mathbf{h}^t} (\mathbf{W}^t)^T \quad \text{and} \quad \frac{\partial \ell}{\partial \mathbf{h}^{t-1}} = \alpha \frac{\partial \ell}{\partial \mathbf{h}'}$$

$$\mathbb{E} \left[\frac{\partial \ell}{\partial h_i^{t-1}} \right] = 0 \qquad \qquad \Rightarrow \qquad \beta = 0$$

$$\text{Var} \left[\frac{\partial \ell}{\partial h_i^{t-1}} \right] = \alpha^2 \text{Var} \left[\frac{\partial \ell}{\partial h_j'} \right] \quad \Rightarrow \quad \alpha = 1$$

Revise Activation Functions

- By the Taylor expansions

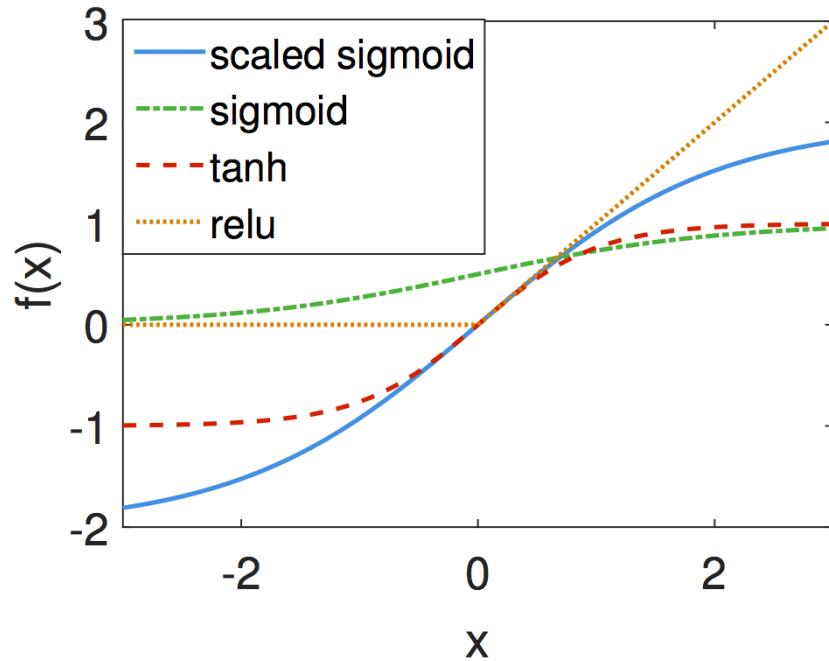
$$\text{sigmoid}(x) = \frac{1}{2} + \frac{x}{4} - \frac{x^3}{48} + O(x^5)$$

$$\tanh(x) = 0 + x - \frac{x^3}{3} + O(x^5)$$

$$\text{relu}(x) = 0 + x \quad \text{for } x \geq 0$$

- “Correct” sigmoid by

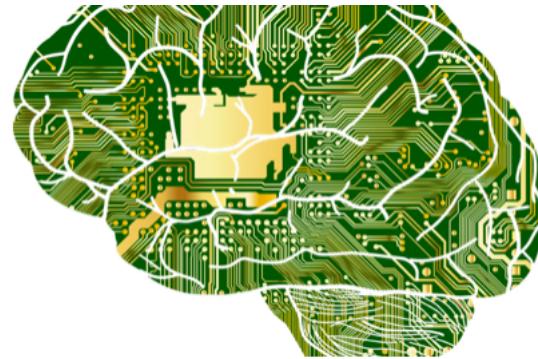
$$4 \times \text{sigmoid}(x) - 2$$



Predicting House Prices on Kaggle

The screenshot shows the Kaggle competition page for "House Prices: Advanced Regression Techniques". At the top, there's a red house icon with a yellow "SOLD" sign. To the right, the title "House Prices: Advanced Regression Techniques" is displayed, along with a brief description: "Predict sales prices and practice feature engineering". Below this, it shows "5,012 teams · Ongoing". A navigation bar follows, with "Overview" underlined in blue, and other tabs: Data, Kernels, Discussion, Leaderboard, Rules, and Team. Under the "Overview" tab, there are several sections: "Description" (highlighted in light blue), "Evaluation", "Frequently Asked Questions", and "Tutorials". To the right, a box titled "Start here if..." contains text about the competition being suitable for data science students who have completed an introductory course and want to expand their skill set before trying a featured competition. At the bottom right, there's a link to the "Competition Description".

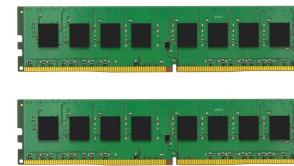
Hardware for Deep Learning



(articleshub360.com)

Your GPU Computer

Intel i7
0.15 TFLOPS

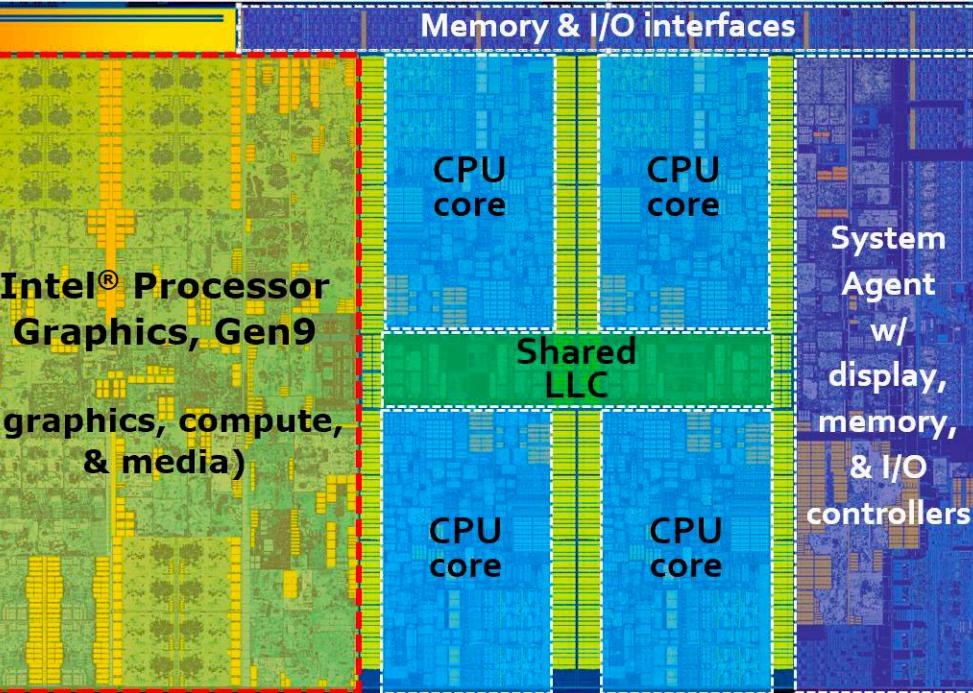


DDR4
32 GB



Nvidia Titan X
12 TFLOPS
16 GB

Intel i7-6700K



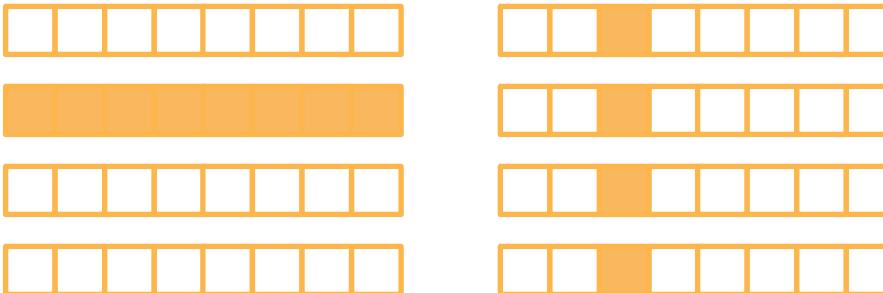
- 4 Physical cores
- Per core
 - 64KB L1 cache
 - 256KB L2 cache
- Shared 8MB L3 cache
- 30 GB/s to main memory

Improve CPU Utilization I

- Before computing $a + b$, need to prepare data first
 - Main memory -> L3 -> L2 -> L1 -> registers
 - L1 cache reference time: 0.5 ns
 - L2 cache reference time 7 ns ($14 \times L1$)
 - Main memory reference time 100ns ($200 \times L1$)
- Improve temporal and spatial memory locality
 - Temporal: reuse data so we keep them on cache
 - Spatial: read data sequential so we can pre-fetch data

Case Study

- For a matrix stored in raw-major, accessing a column is slower than accessing a row
 - CPU reads 64 bytes (cache line) each time
 - CPU “smartly” reads the next cache line ahead when it’s needed



Improve CPU Utilization II

- Server CPUs may have dozens of cores
 - EC2 P3.16xlarge: 2 Intel Xeon CPUs, 32 physical cores in total
- Parallelization to use all cores
 - Hyper-threading may not help because of shared registers

Case Study

- Left is slower than right (Python)

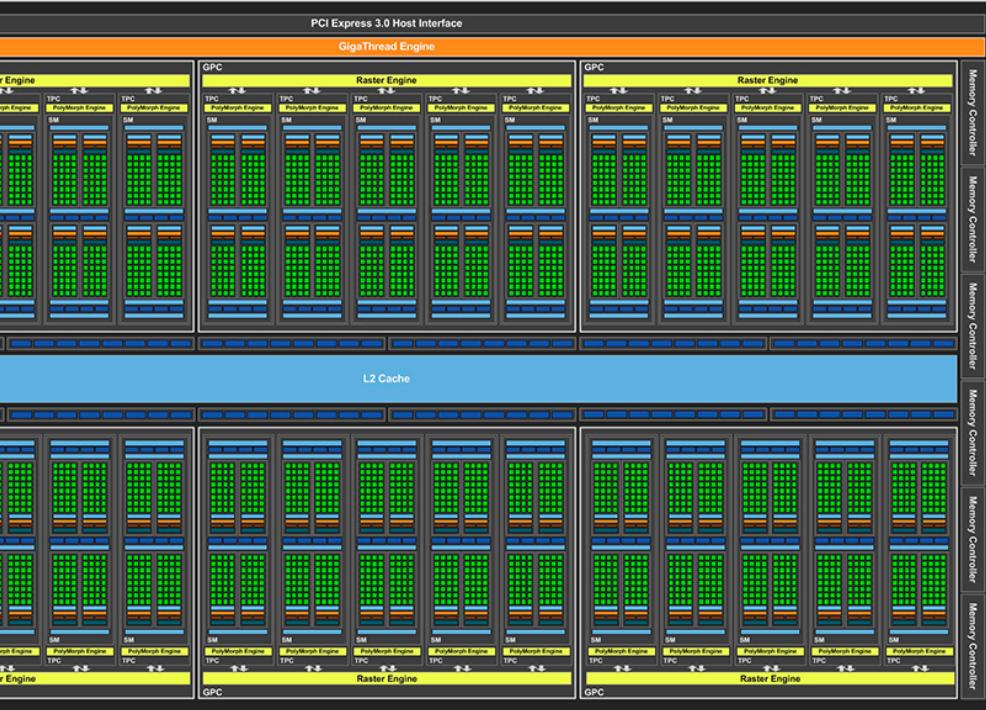
```
for i in range(len(a)):  
    c[i] = a[i] + b[i]
```

```
c = a + b
```

- Left issues n calls, each invoking has certain overhead (e.g. several microsecond)
- Right is easier to be parallelized by a C++ operator

```
#pragma omp for  
for (i=0; i<a.size(); i++) {  
    c[i] = a[i] + b[i]  
}
```

Nvidia Titan X (Pascal)

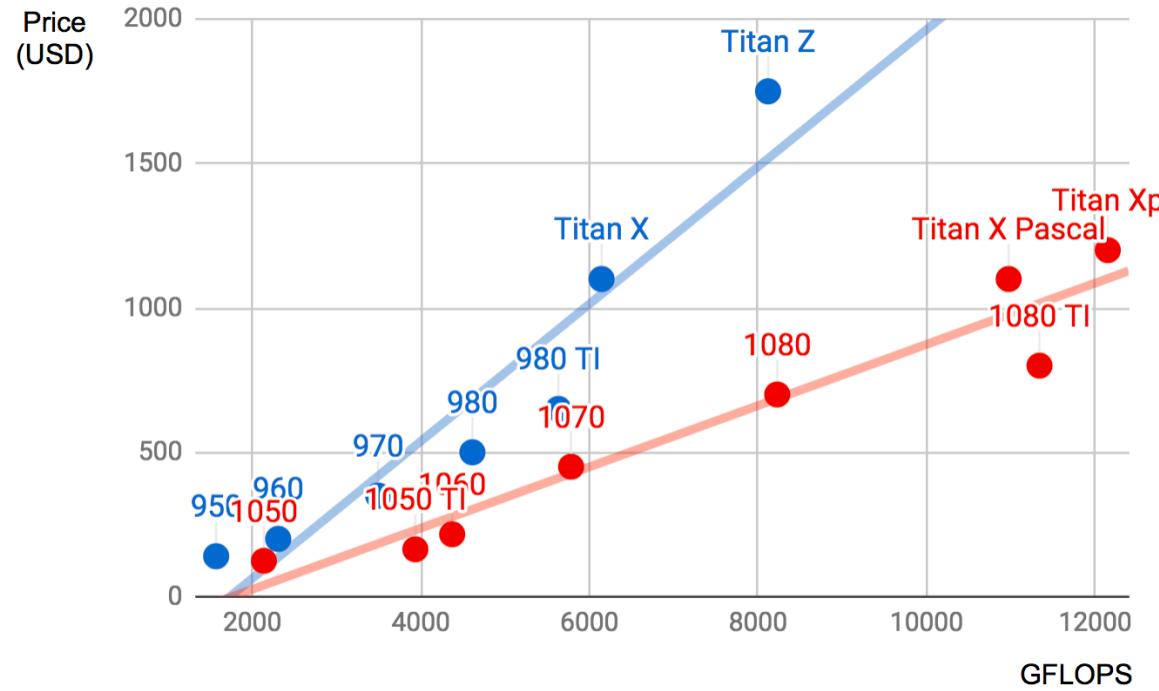


- 2584 cores, each core
 - Multiple arithmetic units
- 3MB L2 cache
- 480 GB/s memory bandwidth

Improve GPU Utilization

- Parallelization
 - Uses thousands of threads
- Memory locality
 - Simple cache architecture and small cache size
- Simple control flows
 - Very limited support
 - Large synchronization

Buy GPUs



- Each series improve the previous one
 - Buy new models
- Price is linear to power in a series

CPU vs GPU



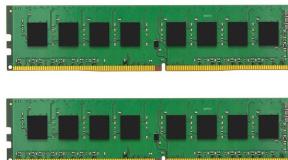
# Cores	6 / 64	2K / 4K
TFLOPS	0.2 / 1	10 / 100
Memory size	32 GB / 1 TB	16 GB / 32 GB
Memory bandwidth	30 GB/s / 100 GB/s	400 GB/s / 1 TB/s
Control flow	Strong	Weak

Typical / High End

CPU/GPU Bandwidth



30 GB/s



PCIe 3.0 16x: 16 GB/s

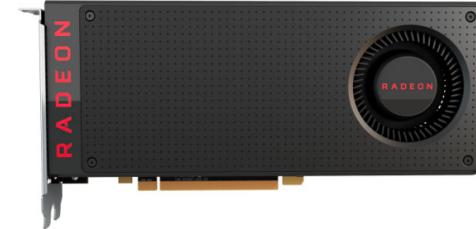


480 GB/s

- Do not copy data between GPU and CPU frequently
 - Limited PCIe bandwidth
 - Synchronization overhead

More CPUs and GPUs

- CPU
 - Desktop/Server: AMD
 - Edge: ARM
- GPU
 - Desktop/Server GPU: AMD, Intel
 - Edge: ARM, Qualcomm



Programming on CPU/GPU

- CPU: C++ or any other high-performance language
 - Mature compilers with performance guaranteed
- GPU
 - CUDA for Nvidia
 - Rich features, mature complier and driver
 - OpenCL for others
 - Quality varies for chip vendors

More Promising Hardware



Qualcomm Snapdragon 845

Snapdragon
X20 LTE modem

Wi-Fi

Hexagon 685 DSP

Qualcomm
Aqstic Audio

System Memory

Adreno 630
Visual Processing
Subsystem

Qualcomm
Spectra 280 ISP

Kryo 385 CPU

Qualcomm
Secure Processing Unit

Touch

PMIC

Digital Signal Processor

- Designed for digital processing algorithms
 - Dot product, Convolution, FFT
- Low power and high performance
 - 5x faster than mobile GPUs, uses less power
- VLIW: Very long instruction word
 - Hundreds multiply-accumulates in a single instruction
- Hard to program and debug
- Compiler toolchain quality varies from chip vendors

Field-Programmable Gate Array (FPGA)

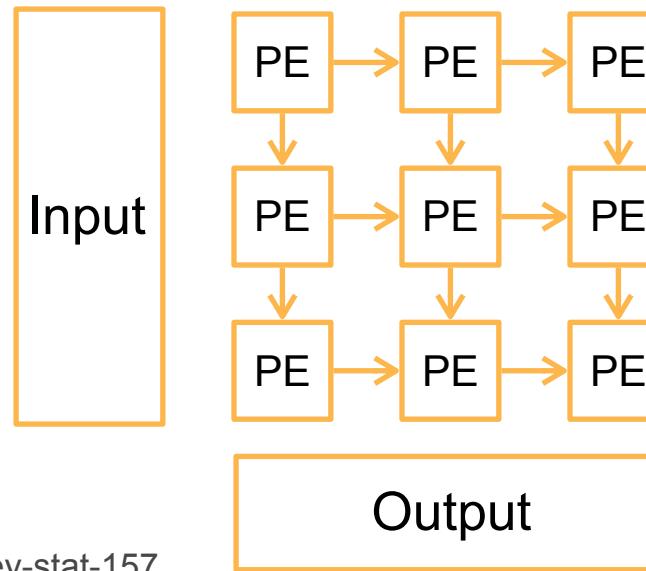
- Contains a large number of programmable logic blocks and reconfigurable interconnects
- Can be configured to perform complex functions
 - Common languages: VHDL and Verilog
- Often has high efficiency than general purpose hardware
- Toolchain qualities vary
- Each “compilation” may take several hours

AI ASIC

- Hot topic in deep learning
 - Every major company is building their chips (Intel, Qualcomm, Google, Amazon, Facebook, ...)
- Google TPU is a landmark chip
 - Matches high-end Nvidia GPU's performance, but 10x-100x cheaper
 - Widely deployed in Google
 - The core is a systolic array

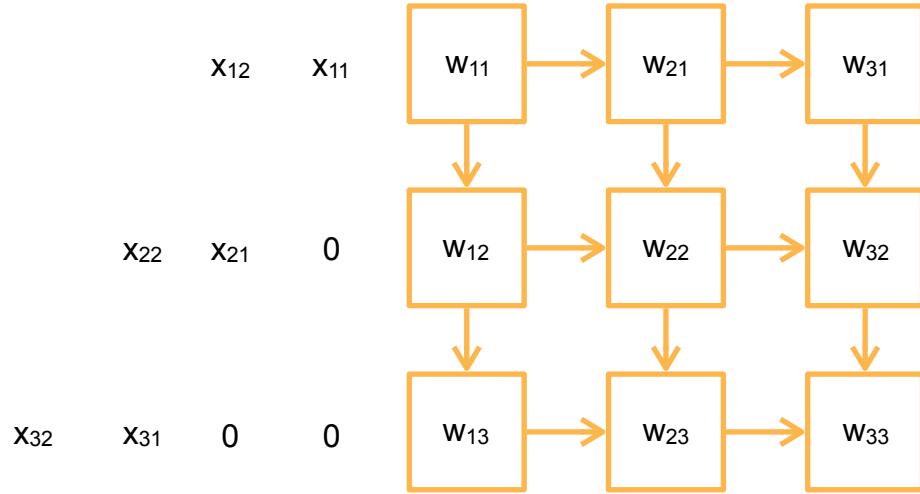
Systolic Array

- An array of processing elements (PE)
- Good at performing matrix-matrix multiplication
- Relative easy/cheaper to build



Matrix Multiplication with Systolic Array

Time 0



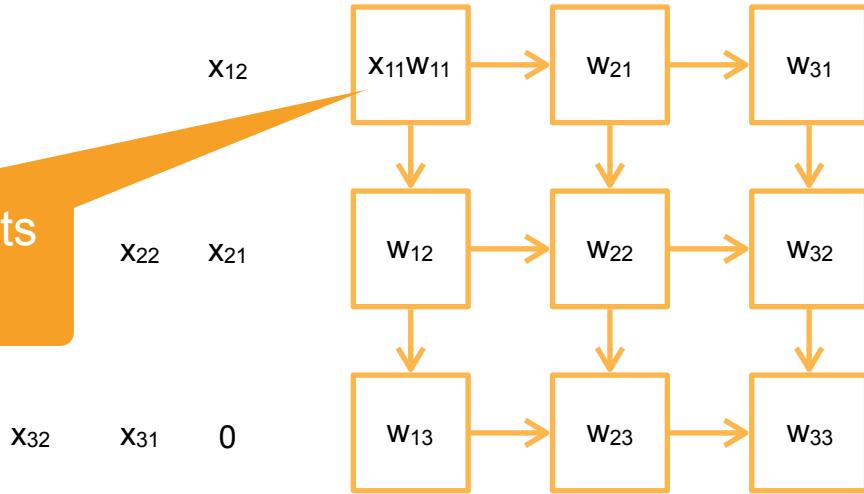
$$\mathbf{Y} = \mathbf{W}\mathbf{X}$$

3x2 3x3 3x2

Matrix Multiplication with Systolic Array

Time 1

Move inputs
to right



$$Y = WX$$

3x2 3x3 3x2

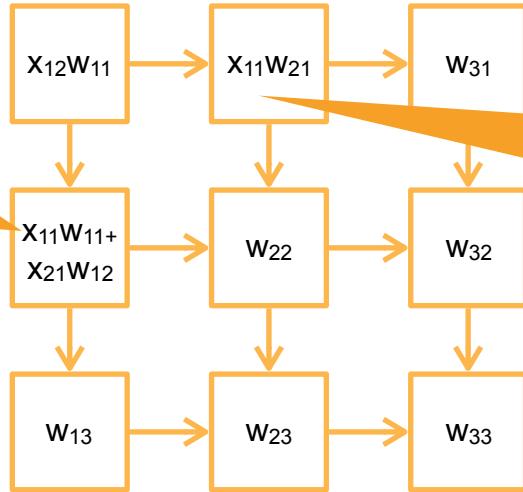
Matrix Multiplication with Systolic Array

Time 2

Move results
to bottom

X₂₂

X₃₂ X₃₁



$$Y = WX$$

3x2 3x3 3x2

Move inputs to
right

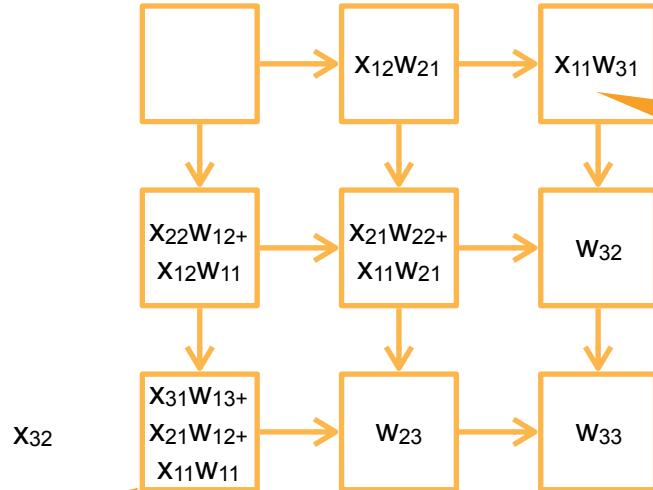
Matrix Multiplication with Systolic Array

Time 3

$$Y = WX$$

3x2 3x3 3x2

Move inputs
to right



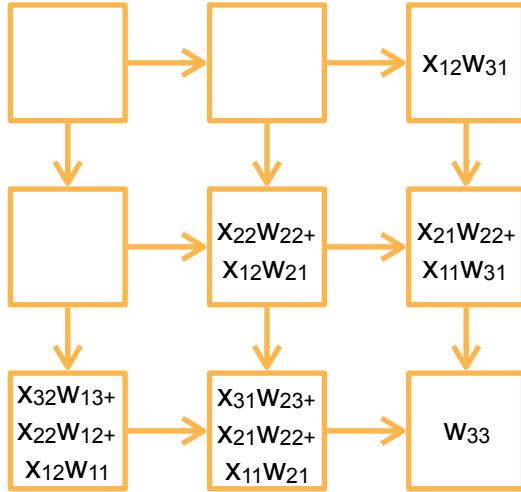
Move results
to bottom

Matrix Multiplication with Systolic Array

Time 4

$$\mathbf{Y} = \mathbf{W}\mathbf{X}$$

3x2 3x3 3x2



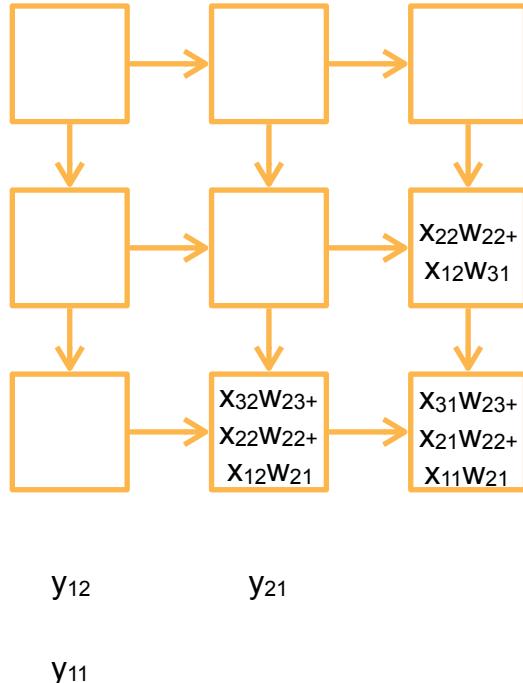
y_{11}

Matrix Multiplication with Systolic Array

Time 5

$$Y = WX$$

3x2 3x3 3x2

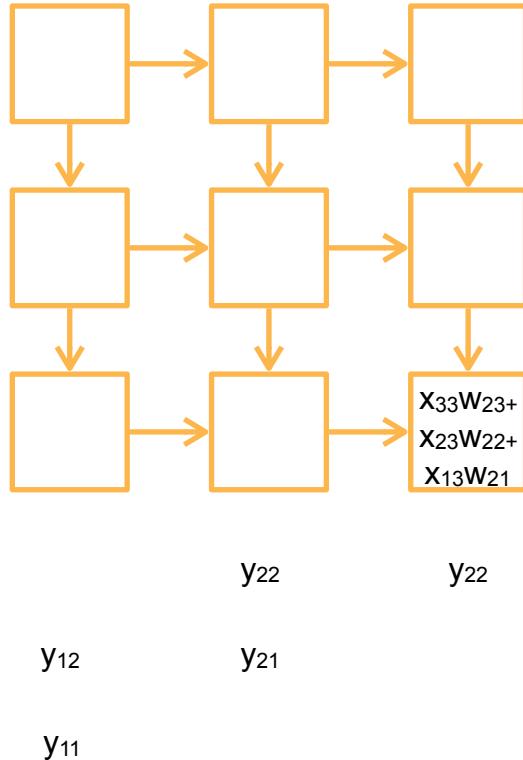


Matrix Multiplication with Systolic Array

Time 6

$$\mathbf{Y} = \mathbf{W}\mathbf{X}$$

3x2 3x3 3x2

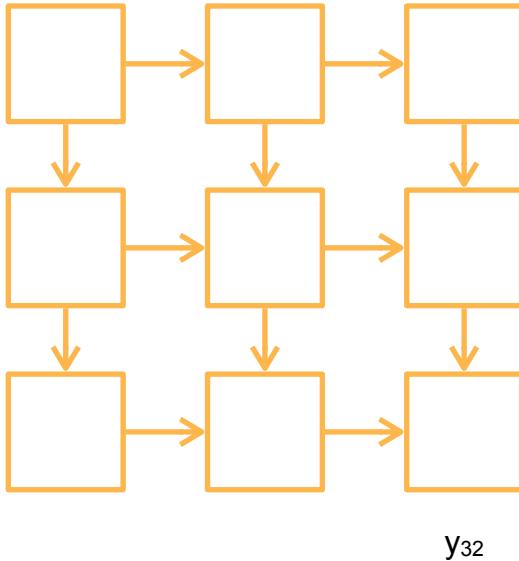


Matrix Multiplication with Systolic Array

Time 7

$$\mathbf{Y} = \mathbf{W}\mathbf{X}$$

3x2 3x3 3x2



Systolic Array

- For general size matrix multiplication, slice and pad inputs to match the SA size
- Batch inputs to reduce the latency cost
- ASIC has other dedicated components for other NN operations, such as sigmoid

Flexibility and Ease of Use



Hexagon 685 DSP



Performance & Power Efficiency