# Copyright Notice

# Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2m} \sum_{j=1}^{n} w_j^2$$

But it makes unacceptably large errors in predictions. What do you try next?

Get more training examples
Try smaller sets of features
Try getting additional features
Try adding polynomial features $(x_1^2, x_2^2, x_1 x_2, etc)$
Try decreasing $\lambda$
Try increasing $\lambda$

# Machine learning diagnostic

Diagnostic: A test that you run to gain insight into what is/isn't working with a learning algorithm, to gain guidance into improving its performance.

Diagnostics can take time to implement but doing so can be a very good use of your time.

Evaluating and
choosing models

Evaluating a model

# Evaluating your model



→ Model fits the training data well but will fail to generalize to new examples not in the training set.

→ $x_1$ = size in feet$^2$
→ $x_2$ = no. of bedrooms
→ $x_3$ = no. of floors
→ $x_4$ = age of home in years

size  $x$

$W_n x^4$

$x = x_1$

$f_{\vec{w},b}(\vec{x}) = w_1 x + w_2 x^2 + \cdots + w_n x^n + b$

just $x$

$f(\vec{x})$

Andrew Ng

# Evaluating your model

Dataset:

| size | price |
|------|-------|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| 1985 | 300 |
| 1534 | 315 |
| 1427 | 199 |
| 1380 | 212 |
| 1494 | 243 |

70%

30%

*training set*

*test set*

$$\left(x^{(1)}, y^{(1)}\right)$$
$$\left(x^{(2)}, y^{(2)}\right)$$
$$\vdots$$
$$\left(x^{(m_{train})}, y^{m_{train}}\right)$$

$m_{train} =$ no. training examples $= 7$

$$\left(x_{test}^{(1)}, y_{test}^{(1)}\right)$$
$$\vdots$$
$$\left(x_{test}^{(m_{test})}, y_{test}^{(m_{test})}\right)$$

$m_{test} =$ no. test examples $= 3$

# Train/test procedure for linear regression (with squared error cost)

Fit parameters by minimizing cost function $J\left(\vec{\mathrm{w}}, b\right)$

$$J\left(\vec{\mathrm{w}}, b\right) = \min_{\vec{\mathrm{w}}, b} \left[ \frac{1}{2m_{train}} \sum_{i=1}^{m_{train}} \left( f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}^{(i)}\right) - y^{(i)} \right)^2 + \frac{\lambda}{2m_{train}} \sum_{j=1}^{n} w_j^2 \right]$$

Compute test error:

$$J_{test}\left(\vec{\mathrm{w}}, b\right) = \frac{1}{2m_{test}} \left[ \sum_{i=1}^{m_{test}} \left( f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}_{test}^{(i)}\right) - y_{test}^{(i)} \right)^2 \right]$$

Compute training error:

$$J_{train}\left(\vec{\mathrm{w}}, b\right) = \frac{1}{2m_{train}} \left[ \sum_{i=1}^{m_{train}} \left( f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}_{train}^{(i)}\right) - y_{train}^{(i)} \right)^2 \right]$$

# Train/test procedure for linear regression (with squared error cost)
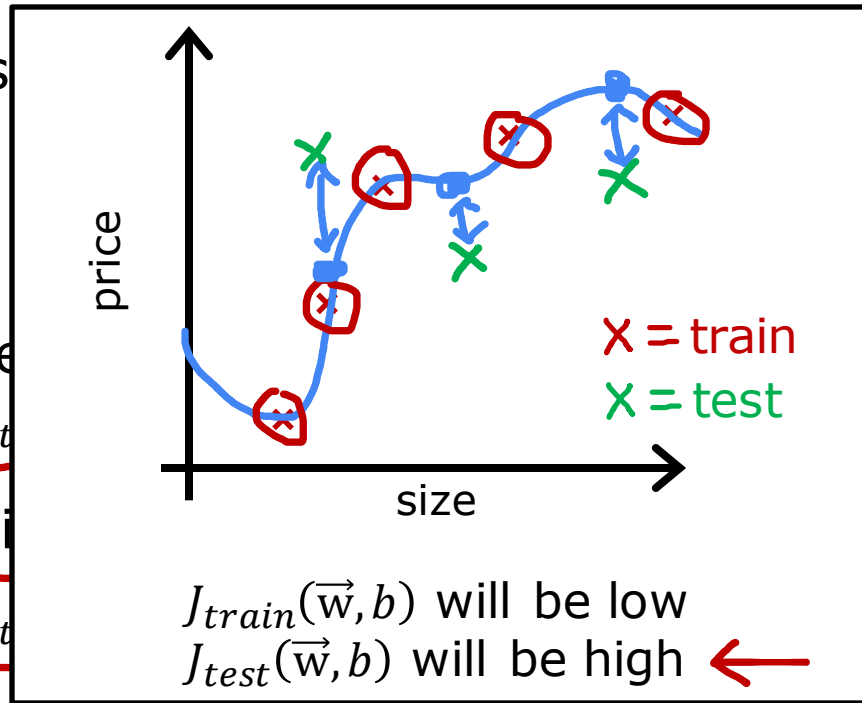
Fit parameters $\quad\quad\quad\quad\quad\quad\quad\quad$ , $b)$

$\longrightarrow J(\vec{w}, b) = \dfrac{min}{\vec{w}, b}$ $\quad\quad\quad\quad\quad\quad$ $\dfrac{}{ain} \displaystyle\sum_{j=1}^{n} w_j^2$ $\Bigg]$ $\longleftarrow$

Compute test e $\quad\quad\quad\quad\quad\quad\quad\quad\quad$ $\Big)^2$ $\Bigg]$ $\quad\quad + \cancel{w^2}$

$\quad\quad\quad J_t \quad\quad\quad\quad\quad\quad\quad\quad$ $\begin{smallmatrix}(i)\\est\end{smallmatrix}$

Compute train $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$

$\quad\quad\quad J_t \quad\quad\quad\quad\quad\quad\quad\quad$ $- y_{train}^{(i)} \Big)^2$ $\Bigg]$



X = train
X = test

$J_{train}(\vec{w}, b)$ will be low
$J_{test}(\vec{w}, b)$ will be high $\longleftarrow$

# Train/test procedure for classification problem

$0/1$

Fit parameters by minimizing $J\left(\vec{\mathrm{w}}, b\right)$ to find $\vec{\mathrm{w}}, b$

E.g.

$$J\left(\vec{\mathrm{w}}, b\right) = -\frac{1}{m}\sum_{i=1}^{m}\left[y^{(i)}\log\left(f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}^{(i)}\right)\right) + (1-y^{(i)})\log\left(1-f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}^{(i)}\right)\right)\right] + \frac{\lambda}{2m}\sum_{j=1}^{n}w_j^2$$

Compute test error:

$$J_{test}\left(\vec{\mathrm{w}}, b\right) = -\frac{1}{m_{test}}\sum_{i=1}^{m_{test}}\left[y_{test}^{(i)}\log\left(f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}_{test}^{(i)}\right)\right) + (1-y_{test}^{(i)})\log\left(1-f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}_{test}^{(i)}\right)\right)\right]$$

Compute train error:

$$J_{train}\left(\vec{\mathrm{w}}, b\right) = -\frac{1}{m_{train}}\sum_{i=1}^{m_{train}}\left[y_{train}^{(i)}\log\left(f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}_{train}^{(i)}\right)\right) + \left(1-y_{train}^{(i)}\right)\log\left(1-f_{\vec{\mathrm{w}},b}\left(\vec{\mathrm{x}}_{train}^{(i)}\right)\right)\right]$$

# Train/test procedure for classification problem $0/1$

Fit parameters by minimizing $J(\overrightarrow{\mathrm{w}}, b)$ to find $\overrightarrow{\mathrm{w}}, b$

E.g.,

$$J(\overrightarrow{\mathrm{w}}, b) = -\frac{1}{m} \sum_{i=1}^{m} \left[ y^{(i)} \log\left( f_{\overrightarrow{\mathrm{w}},b}\left( \cdots \right. \right. \right.$$

Compute test error:

$$J_{test}(\overrightarrow{\mathrm{w}}, b) = -\frac{1}{m_{test}} \sum_{i=1}^{m_{test}} \left[ y \cdots \right.$$

Compute train error:

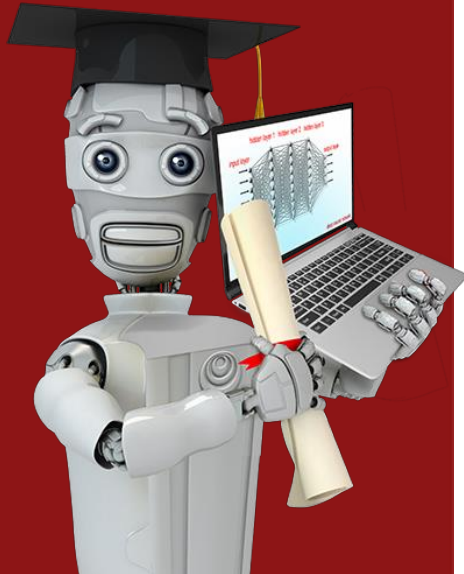$$J_{train}(\overrightarrow{\mathrm{w}}, b) = -\frac{1}{m_{train}} \sum_{i=1}^{m_{train}} \cdots \right) \right]$$

fraction of the test set and the fraction of the train set that the algorithm has misclassified.

$$\hat{y} = \begin{cases} 1 \text{ if } f_{\overrightarrow{\mathrm{w}},b}\left(\overrightarrow{\mathrm{x}}^{(i)}\right) \geq 0.5 \\ 0 \text{ if } f_{\overrightarrow{\mathrm{w}},b}\left(\overrightarrow{\mathrm{x}}^{(i)}\right) < 0.5 \end{cases}$$

count $\hat{y} \neq y$

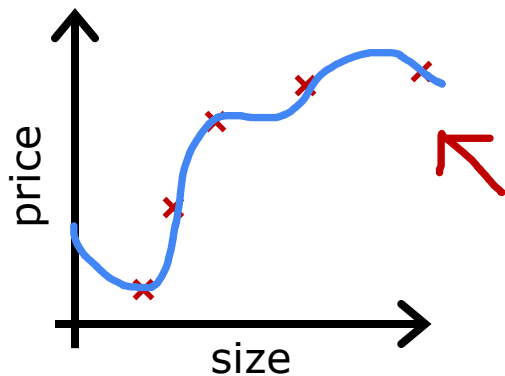$J_{test}(\overrightarrow{\mathrm{w}}, b)$ is the fraction of the test set that has been misclassified.

$J_{train}(\overrightarrow{\mathrm{w}}, b)$ is the fraction of the train set that has been misclassified.

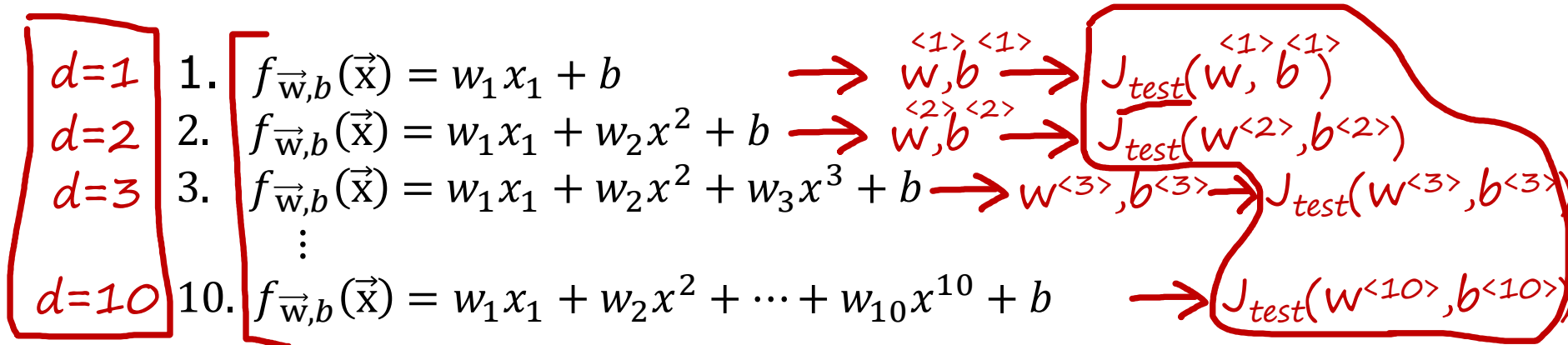# Model selection (choosing a model)



$$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x^2$$
$$+ w_3 x^3 + w_4 x^4 + b$$

Once parameters $\vec{w}, b$ are fit to the training set, the training error $J_{train}(\vec{w}, b)$ is likely lower than the actual generalization error.

$J_{test}(\vec{w}, b)$ is better estimate of how well the model will generalize to new data than $J_{train}(\vec{w}, b)$.
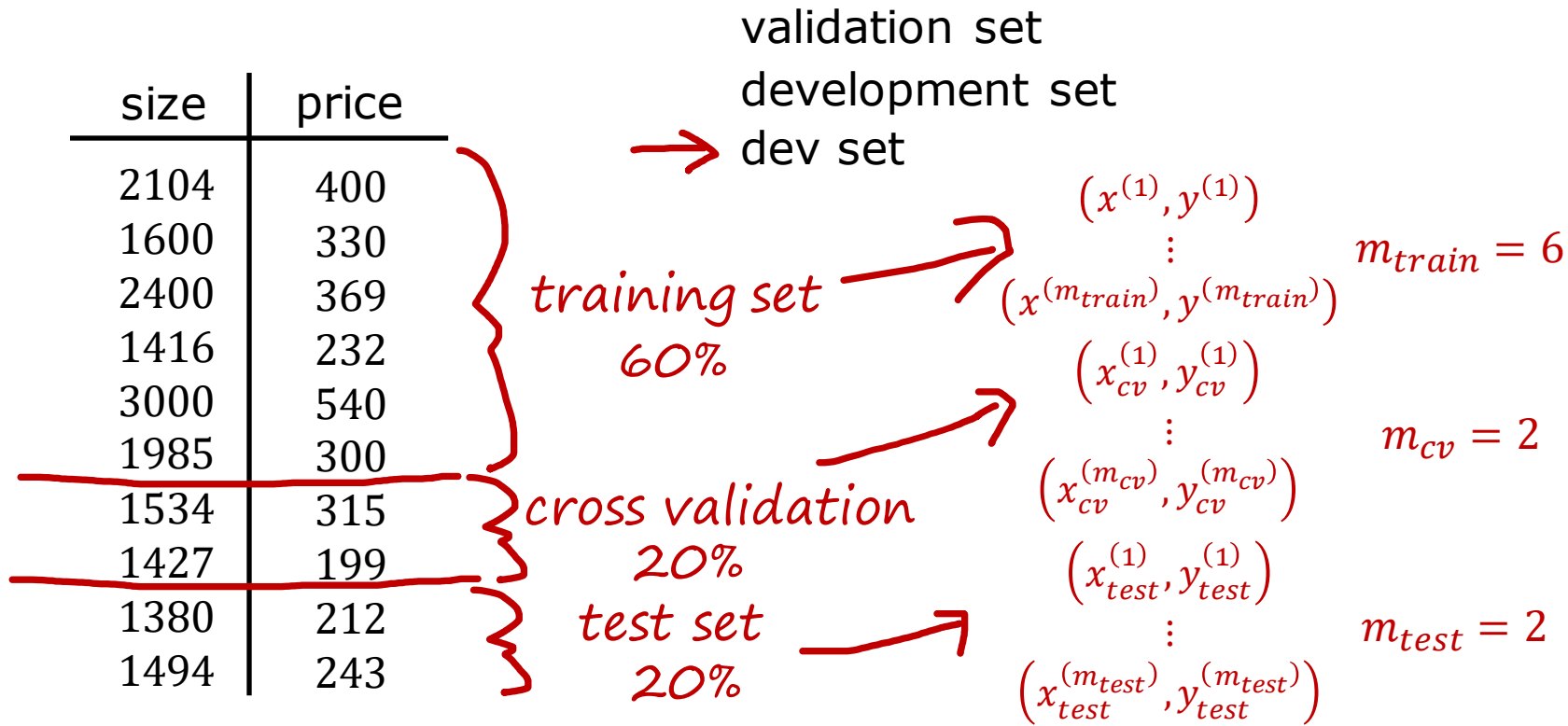
Andrew Ng

# Model selection (choosing a model)

$d=1$   1.   $f_{\vec{w},b}(\vec{x}) = w_1 x_1 + b$    →   $w, b^{<1>}$ → $J_{test}(w^{<1>}, b^{<1>})$

$d=2$   2.   $f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x^2 + b$ → $w, b^{<2>}$ → $J_{test}(w^{<2>}, b^{<2>})$

$d=3$   3.   $f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x^2 + w_3 x^3 + b$ → $w^{<3>}, b^{<3>}$ → $J_{test}(w^{<3>}, b^{<3>})$

      ⋮

$d=10$   10.   $f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x^2 + \cdots + w_{10} x^{10} + b$ → $J_{test}(w^{<10>}, b^{<10>})$

Choose $\boxed{w_1 x_1 + \cdots + w_5 x^5 + b}$   $d=5$     $J_{test}(w^{<5>}, b^{<5>})$

How well does the model perform? Report test set error $J_{test}(w^{<5>}, b^{<5>})$?
The problem is $J_{test}(w^{<5>}, b^{<5>})$ is likely to be an optimistic estimate of
generalization error. Ie: An extra parameter $d$ (degree of polynomial) was
chosen using the test set.     $w, b$

# Training/cross validation/test set

validation set
development set
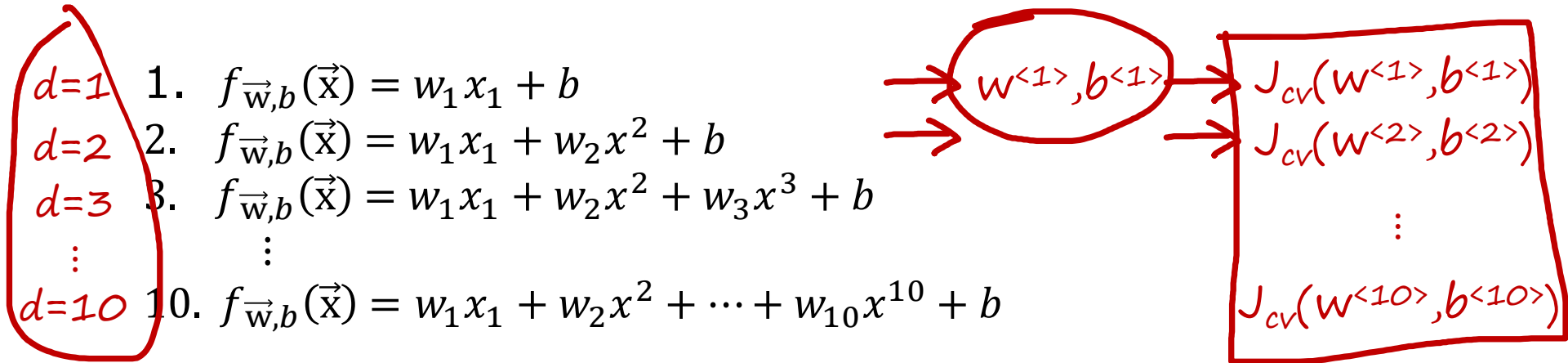dev set

| size | price |
|------|-------|
| 2104 | 400 |
| 1600 | 330 |
| 2400 | 369 |
| 1416 | 232 |
| 3000 | 540 |
| 1985 | 300 |
| 1534 | 315 |
| 1427 | 199 |
| 1380 | 212 |
| 1494 | 243 |

training set
60%

$\left(x^{(1)}, y^{(1)}\right)$
$\vdots$
$\left(x^{(m_{train})}, y^{(m_{train})}\right)$

$m_{train} = 6$

cross validation
20%

$\left(x_{cv}^{(1)}, y_{cv}^{(1)}\right)$
$\vdots$
$\left(x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})}\right)$

$m_{cv} = 2$

test set
20%

$\left(x_{test}^{(1)}, y_{test}^{(1)}\right)$
$\vdots$
$\left(x_{test}^{(m_{test})}, y_{test}^{(m_{test})}\right)$

$m_{test} = 2$

DeepLearning.AI    Stanford ONLINE                              Andrew Ng

# Training/cross validation/test set

Training error:
$$J_{train}(\vec{w}, b) = \frac{1}{2m_{train}}\left[\sum_{i=1}^{m_{train}} \left(f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)}\right)^2\right]$$

Cross validation error:
$$J_{cv}(\vec{w}, b) = \frac{1}{2m_{cv}}\left[\sum_{i=1}^{m_{cv}} \left(f_{\vec{w},b}\left(\vec{x}_{cv}^{(i)}\right) - y_{cv}^{(i)}\right)^2\right]$$
(validation error, dev error)

Test error:
$$J_{test}(\vec{w}, b) = \frac{1}{2m_{test}}\left[\sum_{i=1}^{m_{test}} \left(f_{\vec{w},b}\left(\vec{x}_{test}^{(i)}\right) - y_{test}^{(i)}\right)^2\right]$$

Andrew Ng

# Model selection

$d=1$  1. $f_{\overrightarrow{\mathrm{w}},b}(\vec{\mathrm{x}}) = w_1 x_1 + b$

$d=2$  2. $f_{\overrightarrow{\mathrm{w}},b}(\vec{\mathrm{x}}) = w_1 x_1 + w_2 x^2 + b$

$d=3$  3. $f_{\overrightarrow{\mathrm{w}},b}(\vec{\mathrm{x}}) = w_1 x_1 + w_2 x^2 + w_3 x^3 + b$

$\vdots$

$d=10$  10. $f_{\overrightarrow{\mathrm{w}},b}(\vec{\mathrm{x}}) = w_1 x_1 + w_2 x^2 + \cdots + w_{10} x^{10} + b$

$w^{<1>}, b^{<1>} \longrightarrow J_{cv}(w^{<1>}, b^{<1>})$

$\longrightarrow J_{cv}(w^{<2>}, b^{<2>})$

$\vdots$

$J_{cv}(w^{<10>}, b^{<10>})$

$\longrightarrow$ Pick $w_1 x_1 + \cdots + w_4 x^4 + b$      $(J_{cv}(w^{<4>}, b^{<4>}))$

Estimate generalization error using test the set: $J_{test}(w^{<4>}, b^{<4>})$

# Model selection – choosing a neural network architecture



1.   $w^{(1)}, b^{(1)}$          $J_{cv}(\mathbf{W}^{(1)}, \mathbf{B}^{(1)})$

2.   $w^{(2)}, b^{(2)}$          $J_{cv}(\mathbf{W}^{(2)}, \mathbf{B}^{(2)})$

3.   $w^{(3)}, b^{(3)}$          $J_{cv}(\mathbf{W}^{(3)}, \mathbf{B}^{(3)})$

*Train, CV*

Pick $\mathbf{W}^{(2)}, \mathbf{B}^{(2)}$

Estimate generalization error using the test set: $J_{test}(\mathbf{W}^{(2)}, \mathbf{B}^{(2)})$

# Bias and variance

## Diagnosing bias and variance

# Bias/variance



$$f_{\overrightarrow{\mathbf{w}},b}(x) = w_1 x + b$$

High bias
(underfit)

$$f_{\overrightarrow{\mathbf{w}},b}(x) = w_1 x + w_2 x^2 + b$$

"Just right"

$$f_{\overrightarrow{\mathbf{w}},b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

High variance
(overfit)

$d = 1$    $J_{train}$ is high    $J_{cv}$ is high

$d = 2$    $J_{train}$ is low    $J_{cv}$ is low

$d = 4$    $J_{train}$ is low    $J_{cv}$ is high

# Understanding bias and variance



$J_{cv}(\overrightarrow{w}, b)$

$J_{train}(\overrightarrow{w}, b)$

degree of polynomial

$d$

# Bias and variance

## Regularization and bias/variance

# Linear regression with regularization

Model: $f_{\vec{w},b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2m} \sum_{j=1}^{n} w_j^2$$
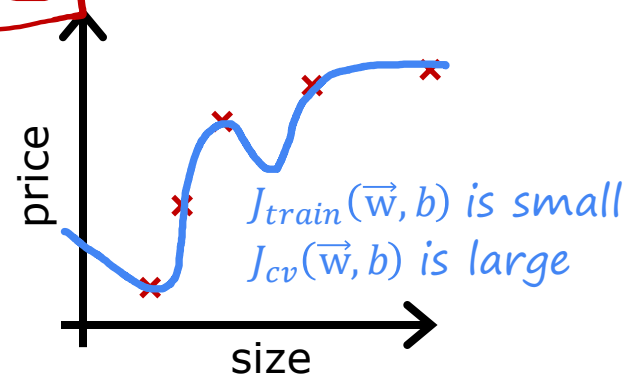


$J_{train}(\vec{w}, b)$ is large

$J_{train}(\vec{w}, b)$ is small
$J_{cv}(\vec{w}, b)$ is small

$J_{train}(\vec{w}, b)$ is small
$J_{cv}(\vec{w}, b)$ is large

Large $\lambda$
High bias (underfit)
$\lambda = 10{,}000$   $w_1 \approx 0, w_2 \approx 0$
$f_{\vec{w},b}(\vec{x}) \approx b$

Intermediate $\lambda$
$\lambda$

Small $\lambda$
High variance (overfit)
$\lambda = 0$

# Choosing the regularization parameter $\lambda$

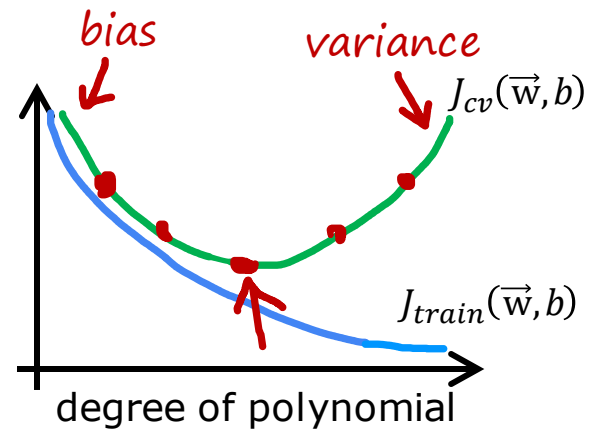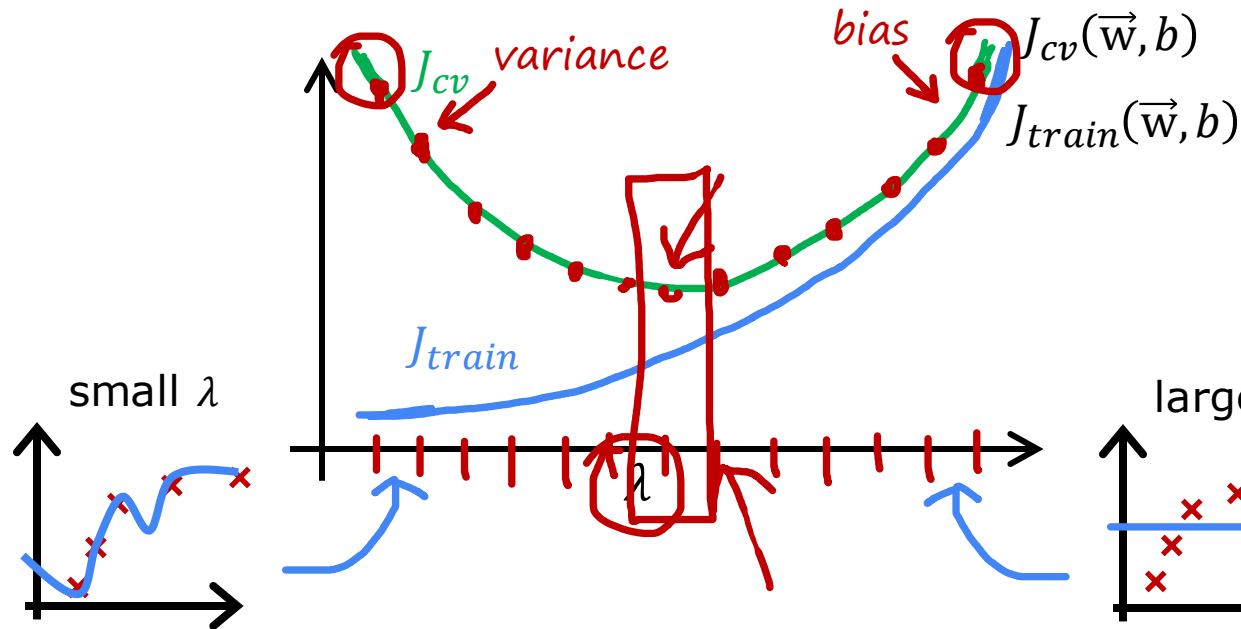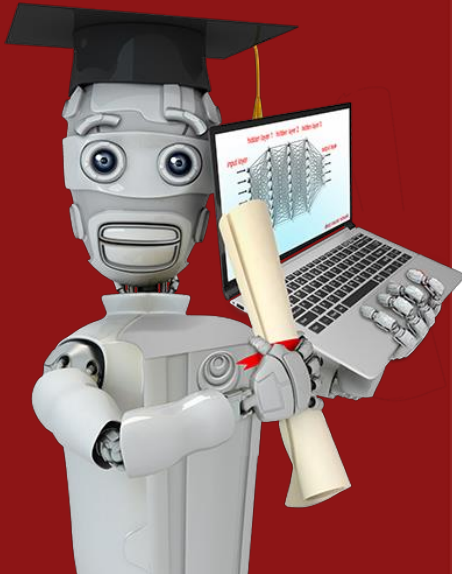Model: $f_{\overrightarrow{w},b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$

1. Try $\lambda = 0$ $\quad\rightarrow\quad \underset{\overrightarrow{w},\,b}{min} J(\overrightarrow{w},b) \quad\rightarrow\quad w^{<1>},b^{<1>} \quad\rightarrow\quad J_{cv}(w^{<1>},b^{<1>})$

2. Try $\lambda = 0.01$ $\quad\rightarrow\quad\quad\quad\quad\quad\rightarrow\quad w^{<2>},b^{<2>} \quad\rightarrow\quad J_{cv}(w^{<2>},b^{<2>})$

3. Try $\lambda = 0.02$ $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\rightarrow\quad J_{cv}(w^{<3>},b^{<3>})$

4. Try $\lambda = 0.04$

5. Try $\lambda = 0.08$ $\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad J_{cv}(w^{<5>},b^{<5>})$

⋮

12. Try $\lambda \approx 10$ $\quad\rightarrow\quad w^{<12>},b^{<12>} \quad\rightarrow\quad J_{cv}(w^{<12>},b^{<12>})$

Pick $w^{<5>},b^{<5>}$

Report test error: $J_{test}(w^{<5>},b^{<5>})$

# Bias and variance as a function of regularization parameter $\lambda$

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^{m} \left(f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)}\right)^2 + \frac{\lambda}{2m} \sum_{j=1}^{n} w_j^2$$



bias

variance

$J_{cv}(\vec{w}, b)$

$J_{train}(\vec{w}, b)$

degree of polynomial

$J_{cv}$ variance

bias $J_{cv}(\vec{w}, b)$

$J_{train}(\vec{w}, b)$

$J_{train}$

$\lambda$

small $\lambda$

large $\lambda$

# Bias and variance

## Establishing a baseline level of performance

# Speech recognition example



Human level performance : $10.6\%$
Training error $J_{train}$ : $10.8\%$
Cross validation error $J_{cv}$ : $14.8\%$

$0.2\%$

$4.0\%$

# Establishing a baseline level of performance

What is the level of error you can reasonably hope to get to?

- Human level performance

- Competing algorithms performance

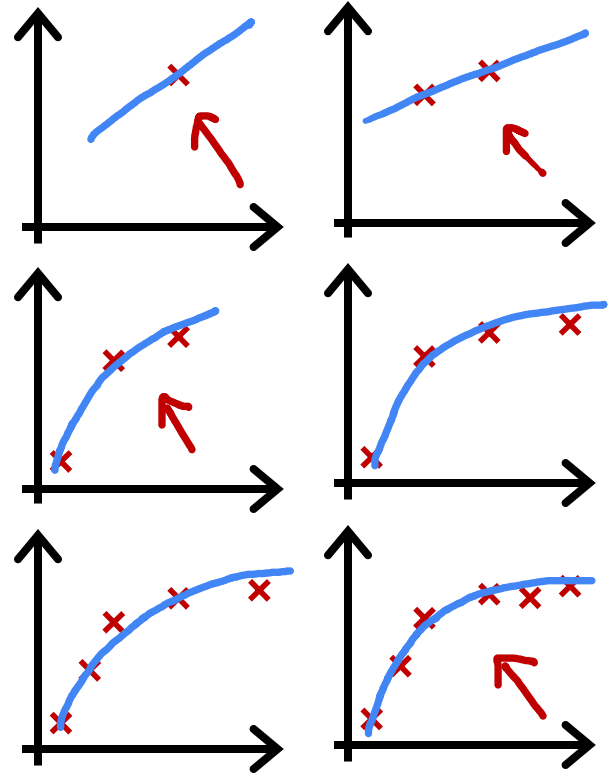- Guess based on experience

# Bias/variance examples

Baseline performance : 10.6% ↕ 0.2%   10.6% ↕ 4.4%   10.6% ↕ 4.4%

Training error （$J_{train}$） : 10.8% ↕ 4.0%   15.0% ↕ 0.5%   15.0% ↕ 4.7%

Cross validation error （$J_{cv}$）: 14.8%    15.5%    19.7%

high variance    high bias    high bias high variance

Bias and variance

Learning curves

# Learning curves

$J_{train}$ = training error

$J_{cv}$ = cross validation error



$$f_{\overrightarrow{w},b}(x) = w_1 x + w_2 x^2 + b$$



error

$J_{cv}(\overrightarrow{w}, b)$

$J_{train}(\overrightarrow{w}, b)$

$m_{train}$ (training set size)

# High bias



error

$J_{cv}(\vec{\mathrm{w}}, b)$

$J_{train}(\vec{\mathrm{w}}, b)$

*human level performance*

$m$ (training set size)

if a learning algorithm suffers from high bias, getting more training data will not (by itself) help much.

$f_{\vec{\mathrm{w}}, b}(x) = w_1 x + b$

# High variance



error

$J_{cv}(\vec{w}, b)$ ←

human level performance

$J_{train}(\vec{w}, b)$

$m$ (training set size)

if a learning algorithm suffers from high variance, getting more training data is likely to help.

$$f_{\vec{w},b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$
(with small $\lambda$)

# Bias and variance

## Deciding what to try next revisited

# Debugging a learning algorithm

You've implemented regularized linear regression on housing prices

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right)^2 + \frac{\lambda}{2m} \sum_{j=1}^{n} w_j^2$$

But it makes unacceptably large errors in predictions. What do you try next?

| | |
|---|---|
| Get more training examples | fixes high variance |
| Try smaller sets of features $\quad x, x^2, x^3, x^4, x^5 \ldots$ | fixes high variance |
| Try getting additional features | fixes high bias |
| Try adding polynomial features $(x_1^2, x_2^2, x_1 x_2, etc)$ | fixes high bias |
| Try decreasing $\lambda$ | fixes high bias |
| Try increasing $\lambda$ | fixes high variance |

# The bias variance tradeoff

$$f_{\overrightarrow{\text{w}},b}(x) = w_1 x + b$$

$$f_{\overrightarrow{\text{w}},b}(x) = w_1 x + w_2 x^2 + b$$

$$f_{\overrightarrow{\text{w}},b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4 + b$$

Simple model

High bias

tradeoff

Complex model

High variance



$J_{cv}(\overrightarrow{\text{w}}, b)$

$J_{train}(\overrightarrow{\text{w}}, b)$

degree of polynomial

# Neural networks and bias variance

Large neural networks are low bias machines



Does it do well on the training set? $J_{train}(\vec{w}, b)$

Yes

Does it do well on the cross validation set? $J_{cv}(\vec{w}, b)$

Yes → Done !

No → Bigger network

No → More data

GPU

# Neural networks and regularization



A large neural network will usually do as well or better than a smaller one so long as regularization is chosen appropriately.

# Neural network regularization

$$J(\mathbf{W}, \mathbf{B}) = \frac{1}{m} \sum_{i=1}^{m} L\big(f(\vec{\mathbf{x}}^{(i)}), y^{(i)}\big) + \frac{\lambda}{2m} \sum_{all\,weights\,\mathbf{W}} (w^2)$$

*b*

## Unregularized MNIST model

```
layer_1 = Dense(units=25, activation="relu")
layer_2 = Dense(units=15, activation="relu")
layer_3 = Dense(units=1, activation="sigmoid")
model = Sequential([layer_1, layer_2, layer_3])
```

## Regularized MNIST model

```
layer_1 = Dense(units=25, activation="relu", kernel_regularizer=L2(0.01))
layer_2 = Dense(units=15, activation="relu", kernel_regularizer=L2(0.01))
layer_3 = Dense(units=1, activation="sigmoid", kernel_regularizer=L2(0.01))
model = Sequential([layer_1, layer_2, layer_3])
```

$\lambda$

# Machine learning development process

## Iterative loop of ML development

# Iterative loop of ML development



Choose architecture
(model, data, etc.)

Train
model

Diagnostics
(bias, variance and
error analysis)

# Spam classification example

From: cheapsales@buystufffromme.com
To: Andrew Ng
Subject: Buy now!

Deal of the week! Buy now!
Rolex w4tchs - $100
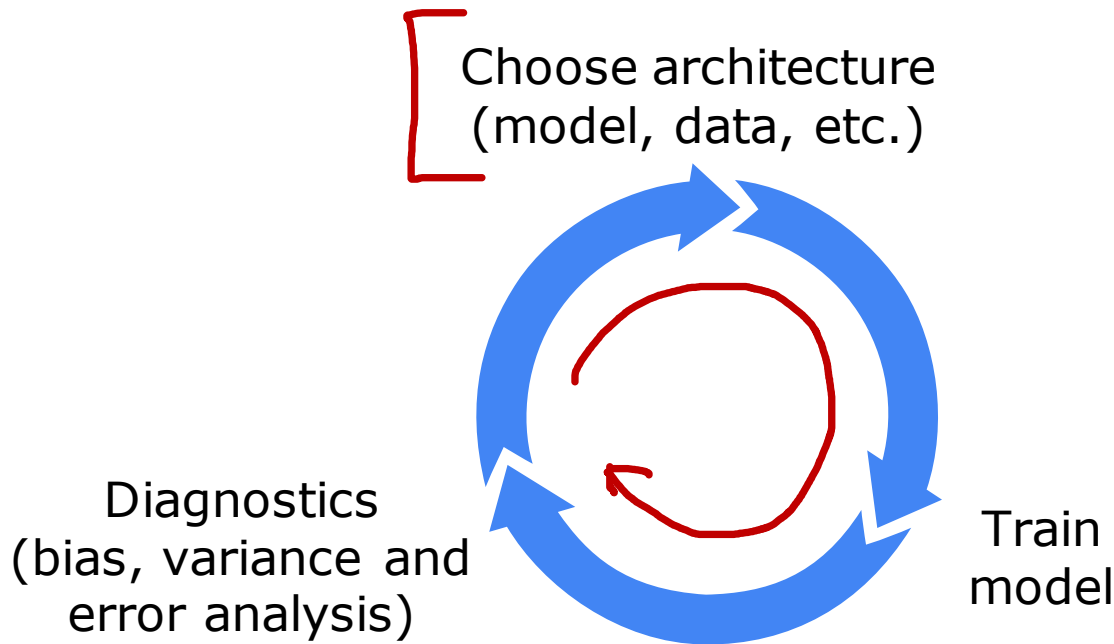Med1cine (any kind) - £50
Also low cost M0rgages
available.

From: Alfred Ng
To: Andrew Ng
Subject: Christmas dates?

Hey Andrew,
Was talking to Mom about plans
for Xmas. When do you get off
work. Meet Dec 22?
Alf

# Building a spam classifier

Supervised learning: $\vec{x}$ = features of email

$y$ = spam (1) or not spam (0)

Features: list the top 10,000 words to compute $x_1, x_2, \cdots, x_{10,000}$

$$\vec{x} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 0 \\ \vdots \end{bmatrix} \begin{matrix} a \\ andrew \\ buy \\ deal \\ discount \\ \vdots \end{matrix}$$

```
From: cheapsales@buystufffromme.com
To: Andrew Ng
Subject: Buy now!

Deal of the week! Buy now!
Rolex w4tchs - $100
Med1cine (any kind) - £50
Also low cost M0rgages
available.
```

# Building a spam classifier

How to try to reduce your spam classifier's error?

- Collect more data. E.g., "Honeypot" project. ⟵
- Develop sophisticated features based on <u>email routing</u> (from <u>email header</u>).
- Define sophisticated features from email body. E.g., should "<u>discounting</u>" and "<u>discount</u>" be treated as the same word.
- Design algorithms to detect misspellings. E.g., w4tches, med1cine, m0rtgage.

# Iterative loop of ML development



Choose architecture
(model, data, etc.)

Train
model

Diagnostics
(bias, variance and
error analysis)

# Error analysis

$m_{cv}$ = ~~500~~ *5000* examples in cross validation set.

Algorithm misclassifies ~~100~~ of them.

*1000*

Manually examine 100 examples and categorize them based on common traits.

Pharma:     *21*     → *more data features*

Deliberate misspellings (w4tches, med1cine): *3*

Unusual email routing:     *7*

Steal passwords (phishing):     *18*     → *more data features*

Spam message in embedded image:     *5*

# Building a spam classifier

How to try to reduce your spam classifier's error?

- Collect more data. E.g., "Honeypot" project.
- Develop sophisticated features based on email routing (from email header).
- Define sophisticated features from email body. E.g., should "discounting" and "discount" be treated as the same word.
- Design algorithms to detect misspellings. E.g., w4tches, med1cine, m0rtgage.

# Iterative loop of ML development



Choose architecture
(model, data, etc.)

Train model

Diagnostics
(bias, variance and error analysis)

Andrew Ng

Machine learning development process

Adding data

# Adding data

→ Add more data of everything. E.g., "Honeypot" project.

→ Add more data of the types where error analysis has indicated it might help.

*Pharma spam*

E.g., Go to unlabeled data and find more examples of Pharma related spam.

( X,Y )

# Data augmentation

Augmentation: modifying an existing training example to create a new training example.

# Data augmentation by introducing distortions

# Data augmentation for speech

Speech recognition example

Original audio (voice search: "What is today's weather?")

+ Noisy background: Crowd

+ Noisy background: Car

+ Audio on bad cellphone connection

# Data augmentation by introducing distortions

Distortion introduced should be representation of the type of noise/distortions in the test set.



Audio:
Background noise,
bad cellphone connection

Usually does not help to add purely random/meaningless noise to your data.



$x_i$ =intensity (brightness) of pixel $i$
$x_i \leftarrow x_i + $ random noise

[Adam Coates and Tao Wang]

# Data synthesis

Synthesis: using artificial data inputs to create a new training example.

# Artificial data synthesis for photo OCR



[http://www.publicdomainpictures.net/view-image.php?image=5745&picture=times-square]

# Artificial data synthesis for photo OCR



Real data

Abcdefg
Abcdefg
Abcdefg
Abcdefg
Abcdefg
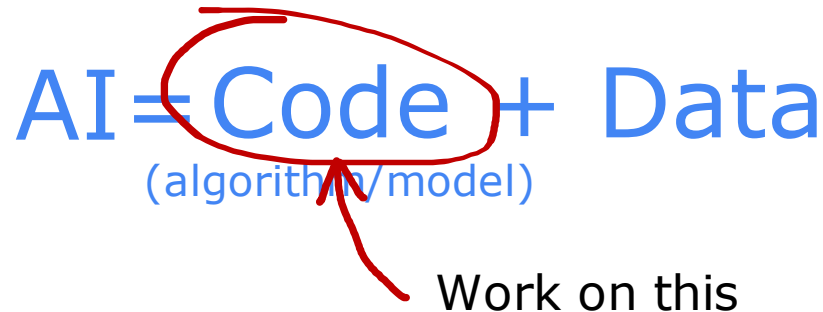
# Artificial data synthesis for photo OCR



Real data



Synthetic data

[Adam Coates and Tao Wang]

# Engineering the data used by your system

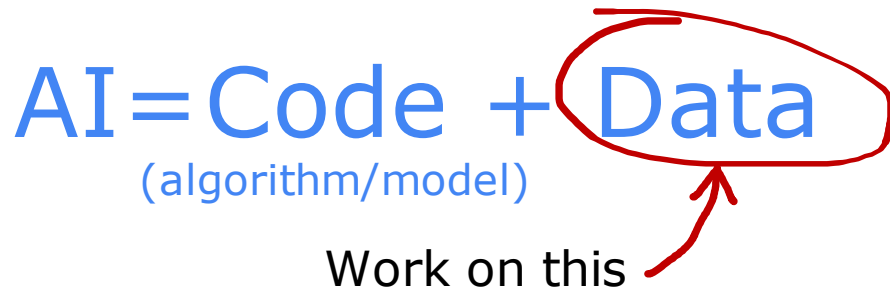Conventional model-centric approach:
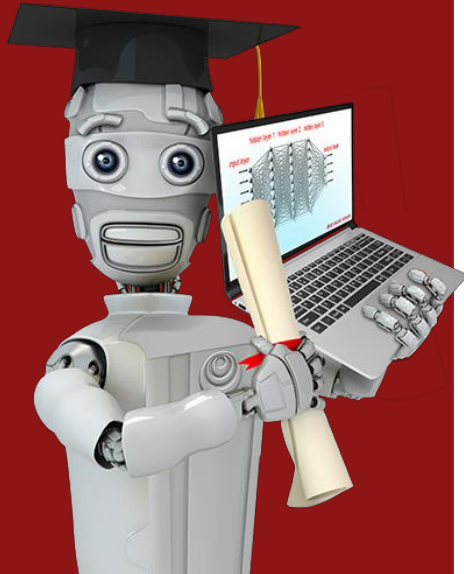
AI = ⟨Code⟩ + Data
(algorithm/model)

Work on this

Data-centric approach:
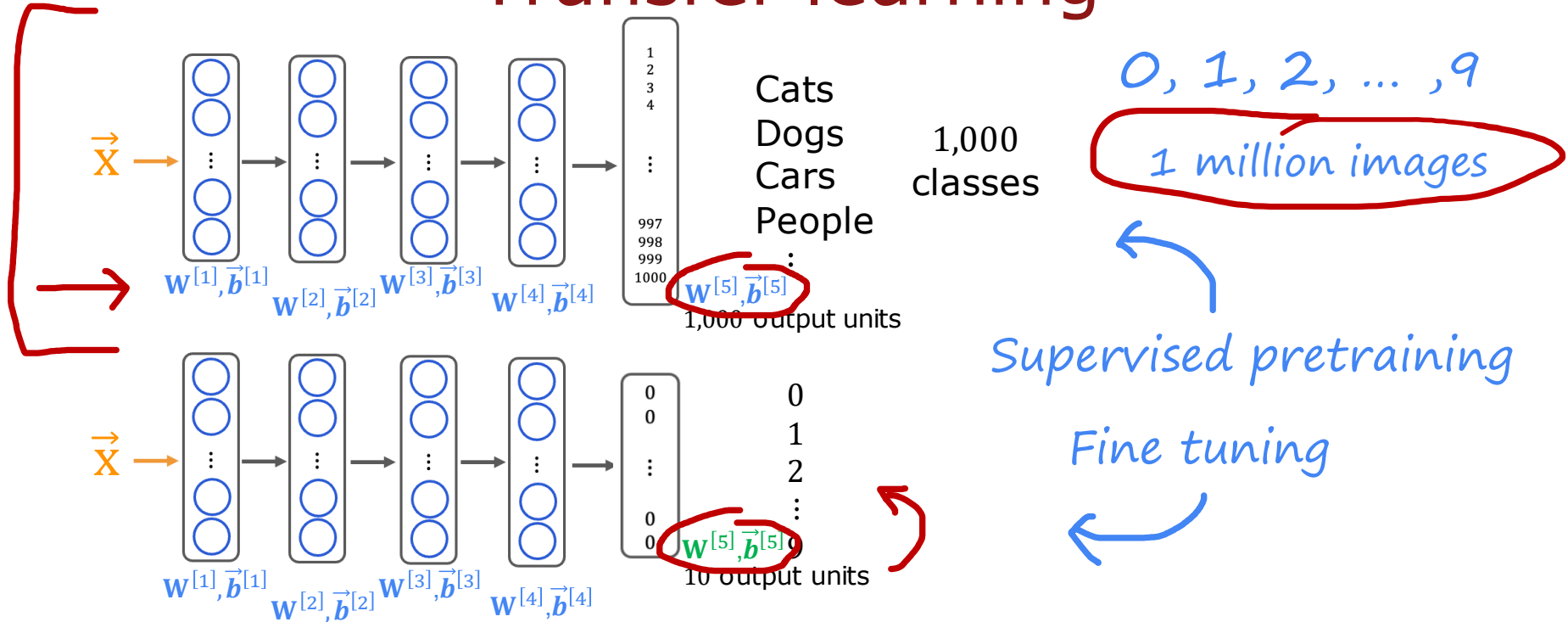
AI = Code + ⟨Data⟩
(algorithm/model)

Work on this

# Transfer learning



**Cats**
**Dogs**
**Cars**
**People**

1,000 classes

$W^{[1]}, \vec{b}^{[1]}$  $W^{[2]}, \vec{b}^{[2]}$  $W^{[3]}, \vec{b}^{[3]}$  $W^{[4]}, \vec{b}^{[4]}$  $W^{[5]}, \vec{b}^{[5]}$

1,000 output units

$0, 1, 2, \dots , 9$

1 million images

Supervised pretraining

Fine tuning

0
1
2
⋮
0

$W^{[1]}, \vec{b}^{[1]}$  $W^{[2]}, \vec{b}^{[2]}$  $W^{[3]}, \vec{b}^{[3]}$  $W^{[4]}, \vec{b}^{[4]}$  $W^{[5]}, \vec{b}^{[5]}$

10 output units
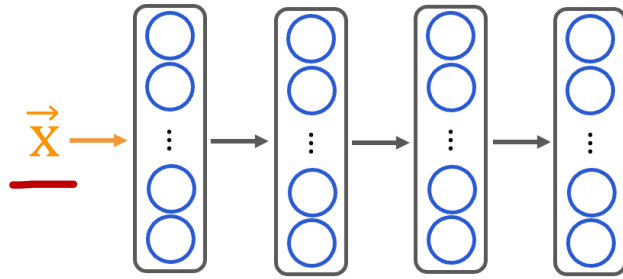
Option 1: only train output layers parameters.
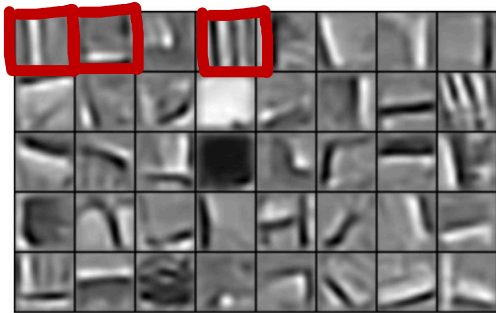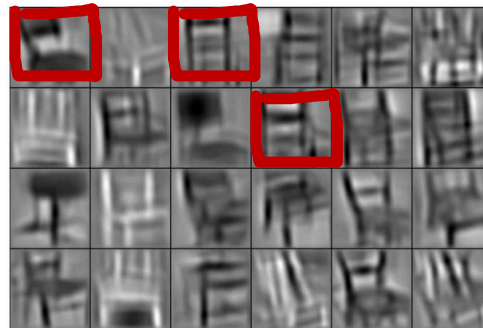
Option 2: train all parameters.

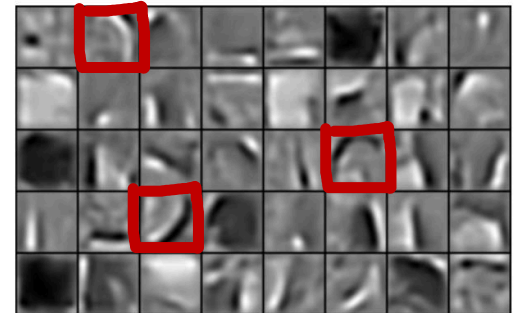# Why does transfer learning work?

use the same input type

detects edges

detects corners

detects curves/basic shapes

Edges

Corners

Curves / basic shapes

DeepLearning.AI    Stanford ONLINE

Andrew Ng

# Transfer learning summary

1. Download neural network parameters pretrained on a large dataset with same input type (e.g., images, audio, text) as your application (or train your own).
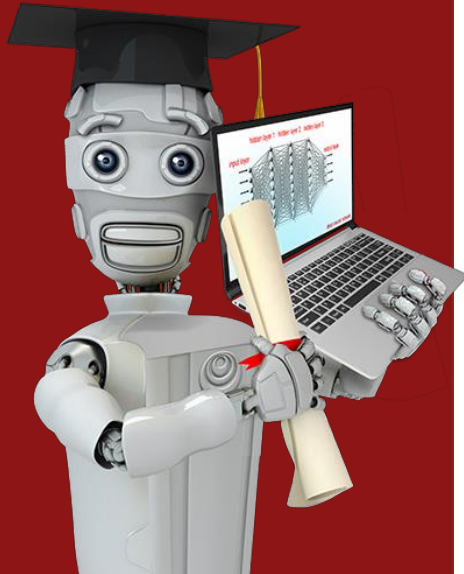
1M

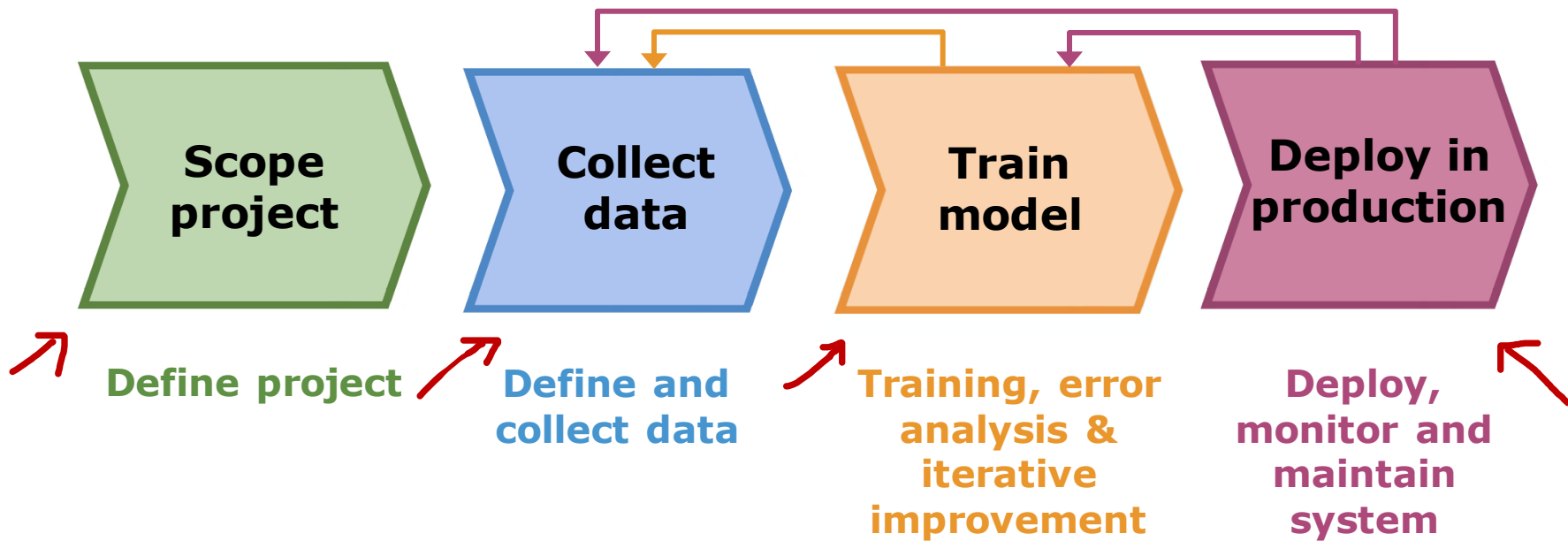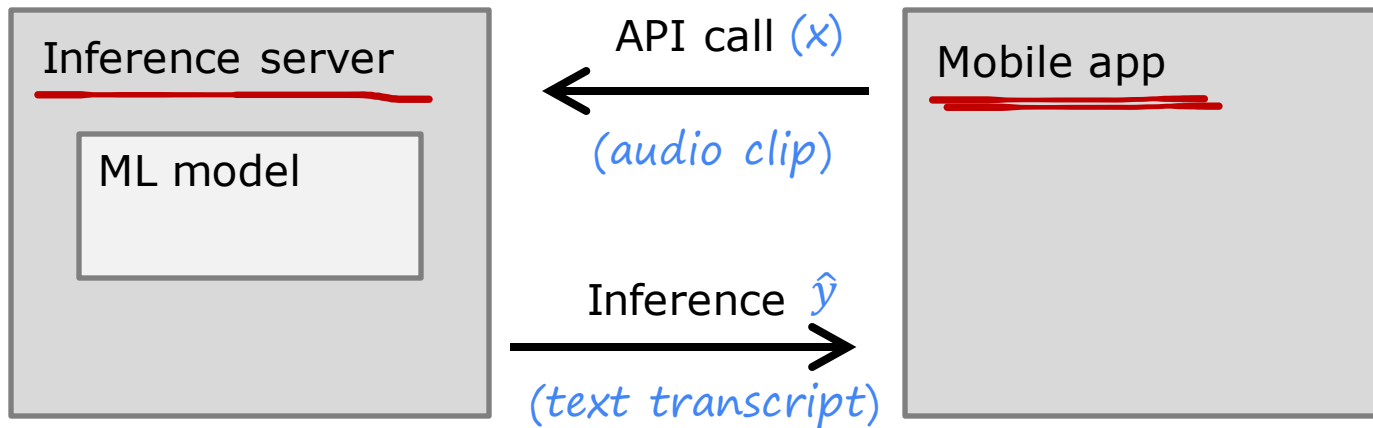2. Further train (fine tune) the network on your own data.

1000

50

# Full cycle of a machine learning project

# Deployment



Inference server

ML model

API call ($x$)

(audio clip)

Mobile app

Inference $\hat{y}$

(text transcript)

Software engineering may be needed for:

Ensure reliable and efficient predictions
Scaling
Logging
System monitoring
Model updates

MLOps
machine learning operations

# Bias

Hiring tool that discriminates against women.

Facial recognition system matching dark skinned individuals to criminal mugshots.

Biased bank loan approvals.

Toxic effect of reinforcing negative stereotypes.

# Adverse use cases

Deepfakes

Spreading toxic/incendiary speech through optimizing for engagement.

Generating fake content for commercial or political purposes.

Using ML to build harmful products, commit fraud etc.
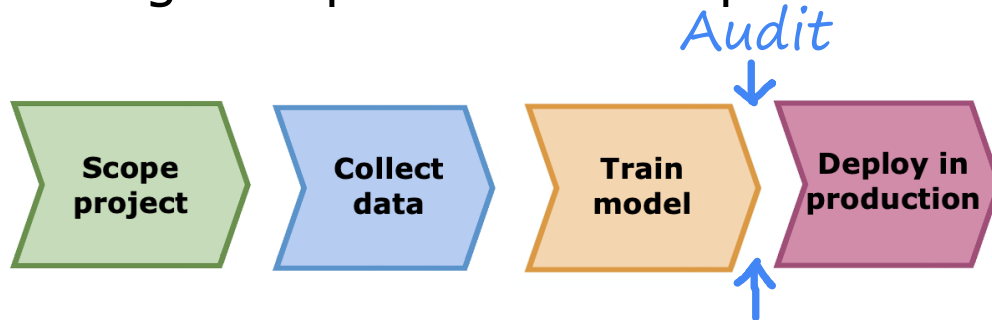         Spam vs anti-spam : fraud vs anti-fraud.

Andrew Ng

# Guidelines

Get a diverse team to brainstorm things that might go wrong, with emphasis on possible harm to vulnerable groups.

Carry out literature search on standards/guidelines for your industry.

Audit systems against possible harm prior to deployment.



Develop mitigation plan (if applicable), and after deployment, monitor for possible harm.

# Skewed datasets (optional)

## Error metrics for skewed datasets

# Rare disease classification example

Train classifier $f_{\vec{\mathbf{w}}, b}(\vec{\mathbf{x}})$     ($y = 1$ if disease present,
$\qquad\qquad\qquad\qquad\qquad\quad y = 0$ otherwise)

Find that you've got 1% error on test set
(99% correct diagnoses)

Only 0.5% of patients have the disease

print("y=0")     99.5% accuracy, 0.5% error

1%

1.2%

# Precision/recall

$y = 1$ in presence of rare class we want to detect.

Actual Class

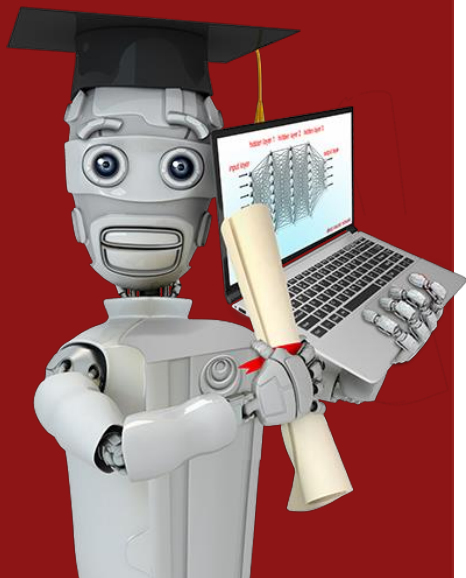|  | 1 | 0 |
|---|---|---|
| Predicted Class 1 | True positive 15 | False positive 5 |
| 0 | False negative 10 | True negative 70 |

↓ 25 ↓ 75

`print("y=0")`

**Precision:**
(of all patients where we predicted $y = 1$, what fraction actually have the rare disease?)

$$\frac{\text{True positives}}{\#\text{predicted positive}} = \frac{\text{True positives}}{\text{True pos} + \text{False pos}} = \frac{15}{15+5} = 0.75$$

**Recall:** ←
(of all patients that actually have the rare disease, what fraction did we correctly detect as having it?)

$$\frac{\text{True positives}}{\#\text{actual positive}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}} = \frac{15}{15+10} = 0.6$$

# Skewed datasets (optional)

## Trading off precision and recall

# Trading off precision and recall

Logistic regression: $0 < f_{\vec{w},b}(\vec{x}) < 1$

Predict 1 if $f_{\vec{w},b}(\vec{x}) \geq 0.5$ ~~0.7~~ ~~0.9~~ 0.3

Predict 0 if $f_{\vec{w},b}(\vec{x}) < 0.5$ ~~0.7~~ ~~0.9~~ 0.3

$$\text{precision} = \frac{\text{true positives}}{\text{total predicted positive}}$$

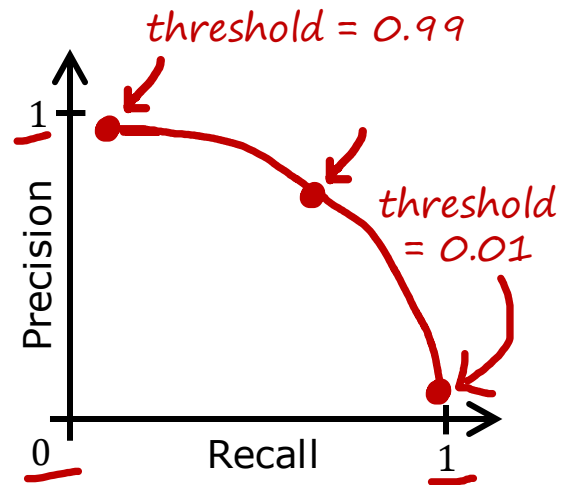$$\text{recall} = \frac{\text{true positives}}{\text{total actual positive}}$$

Suppose we want to predict $y = 1$ (rare disease) only if very confident.

higher precision, lower recall.

Suppose we want to avoid missing too many cases of rare disease (when in doubt predict $y = 1$)

lower precision, higher recall.

More generally predict 1 if: $f_{\vec{w},b}(\vec{x}) \geq$ threshold.



threshold = 0.99

threshold = 0.01

Precision

Recall

# F1 score

How to compare precision/recall numbers?

|  | Precision (P) | Recall (R) | Average | $F_1$ score |
|---|---|---|---|---|
| Algorithm 1 | 0.5 ↔ | 0.4 | ~~0.45~~ | 0.444 ← |
| Algorithm 2 | 0.7 | 0.1 | ~~0.4~~ | 0.175 |
| Algorithm 3 | 0.02 | 1.0 | ~~0.501~~ | 0.0392 |

`print("y=1")`

$$\text{Average} = \frac{P+R}{2}$$

$$\text{F1 score} = \frac{1}{\frac{1}{2}\left(\frac{1}{P}+\frac{1}{R}\right)} = 2\frac{PR}{P+R}$$