# Linear Regression

PROF. SAMEER SINGH

FALL 2017

CS 273A: Machine Learning

# Machine Learning

Bayes Error Wrapup

Gaussian Bayes Models

Linear Regression: Definition and Cost

Gradient Descent Algorithms

# Measuring errors

Confusion matrix

Can extend to more classes

$$accracy = \frac{t_p + t_n}{t_p + t_n + f_p + f_n}$$

$$err = 1 - accuracy$$

|  | Predict: 0 | Predict 1 |
|---|---|---|
| Y=0 | 380 | 5 |
| Y=1 | 338 | 3 |

*True* (vertical label, left of table)

Annotations: $t_n$ → 380, $f_p$ → 5, $f_n$ → 338, $t_p$ → 3

True positive rate: #(y=1 , ŷ=1) / #(y=1)    -- "sensitivity"
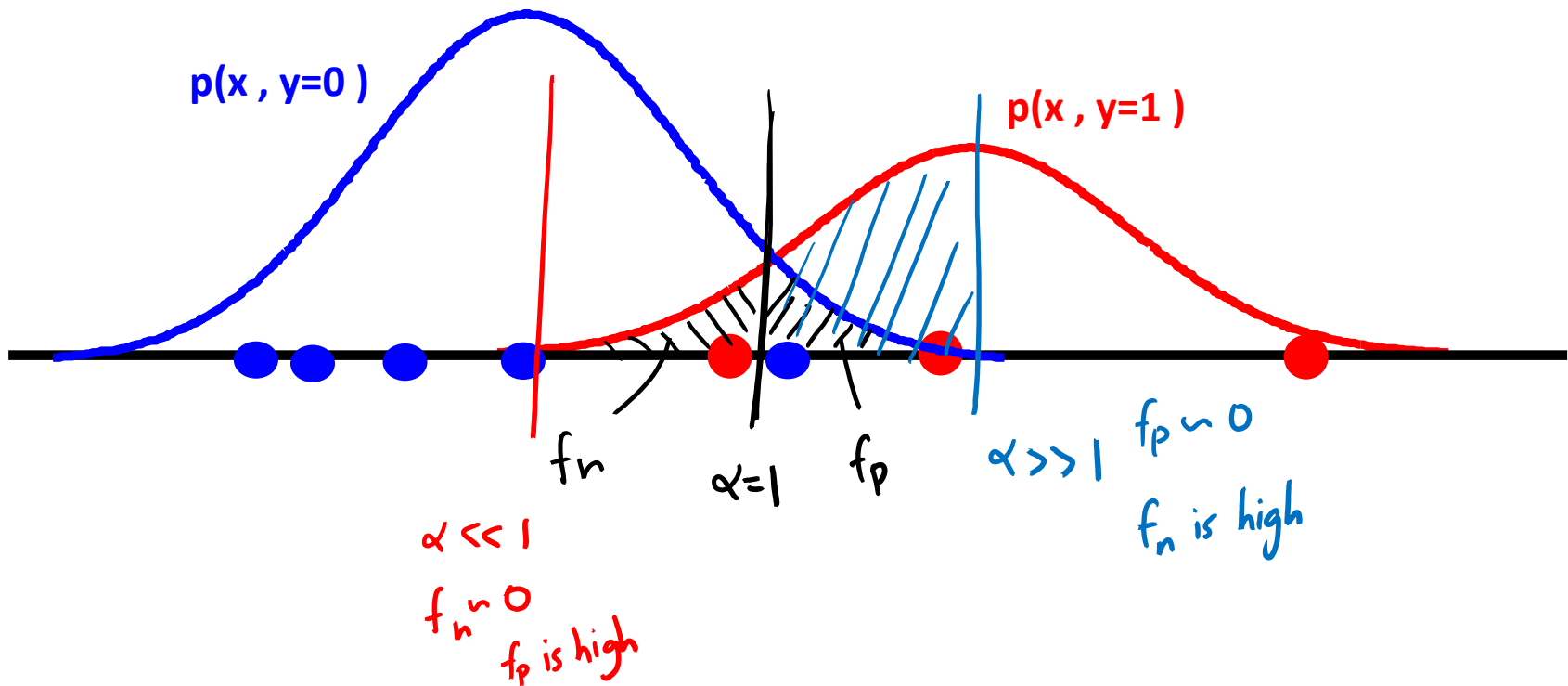
False negative rate:  #(y=1 , ŷ=0) / #(y=1)

False positive rate:  #(y=0 , ŷ=1) / #(y=0)

True negative rate:  #(y=0 , ŷ=0) / #(y=0)    -- "specificity"

# Decision Surfaces

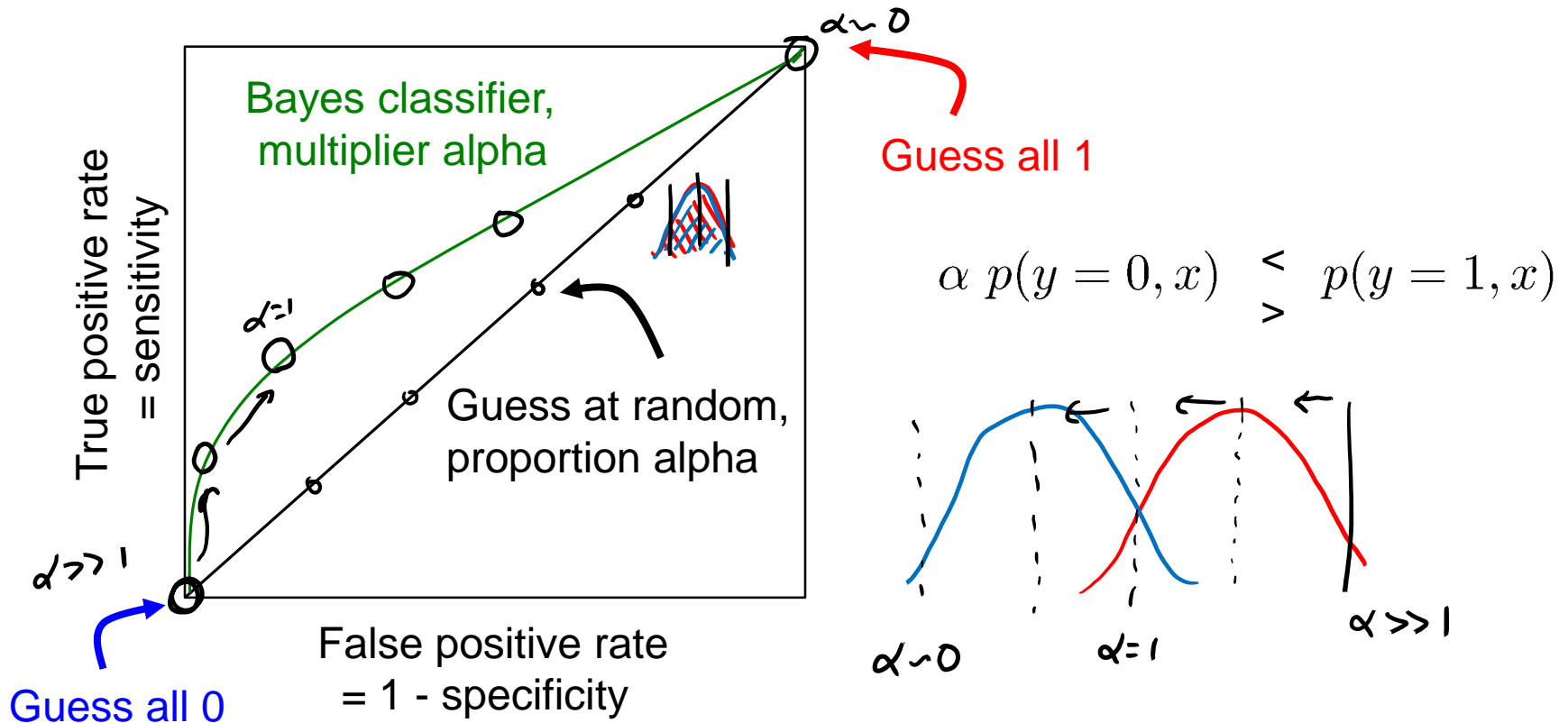Add multiplier alpha:

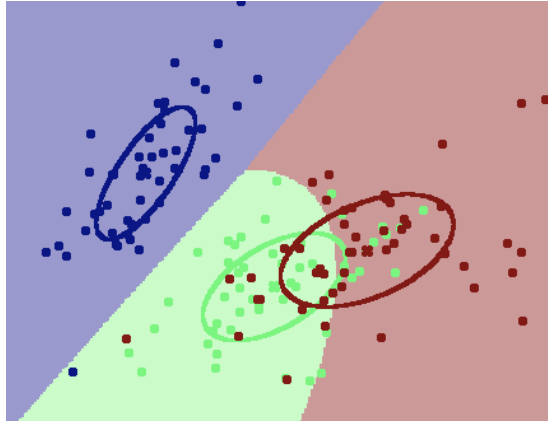$$\alpha \; p(y = 0, x) \begin{array}{c} < \\ > \end{array} p(y = 1, x)$$
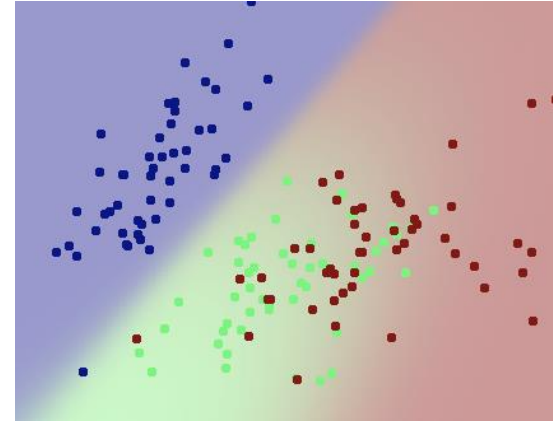


p(x , y=0 )

p(x , y=1 )

fn

$\alpha = 1$

fp

$\alpha >> 1$   fp ~ 0
fn is high

$\alpha << 1$
fn ~ 0
   fp is high

# ROC Curves

Characterize performance as we vary the decision threshold?



$$\alpha \, p(y = 0, x) \quad \substack{< \\ >} \quad p(y = 1, x)$$

# Probabilistic vs. Discriminative learning
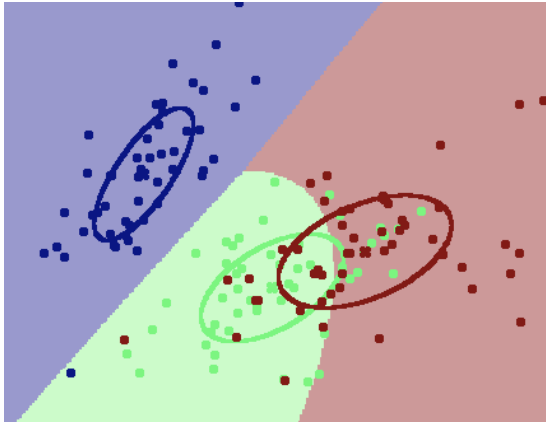


"Discriminative" learning:
Output prediction ŷ(x)



"Probabilistic" learning:
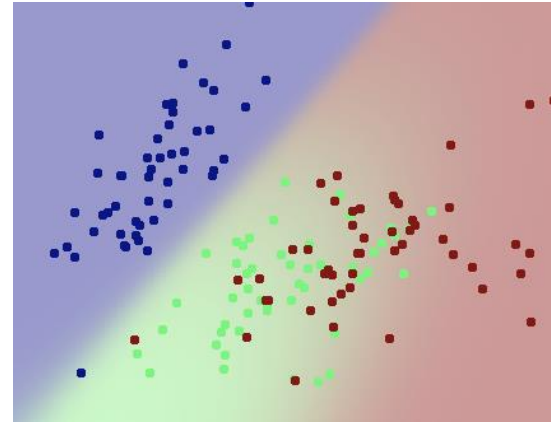Output probability p(y|x)
*(expresses confidence in outcomes)*

"Probabilistic" learning
- Conditional models just explain y:  p(y|x)
- Generative models also explain x: p(x,y)
  - Often a component of unsupervised or semi-supervised learning
- Bayes and Naïve Bayes classifiers are generative models

# Probabilistic vs. Discriminative learning



"Discriminative" learning:
Output prediction ŷ(x)



"Probabilistic" learning:
Output probability p(y|x)
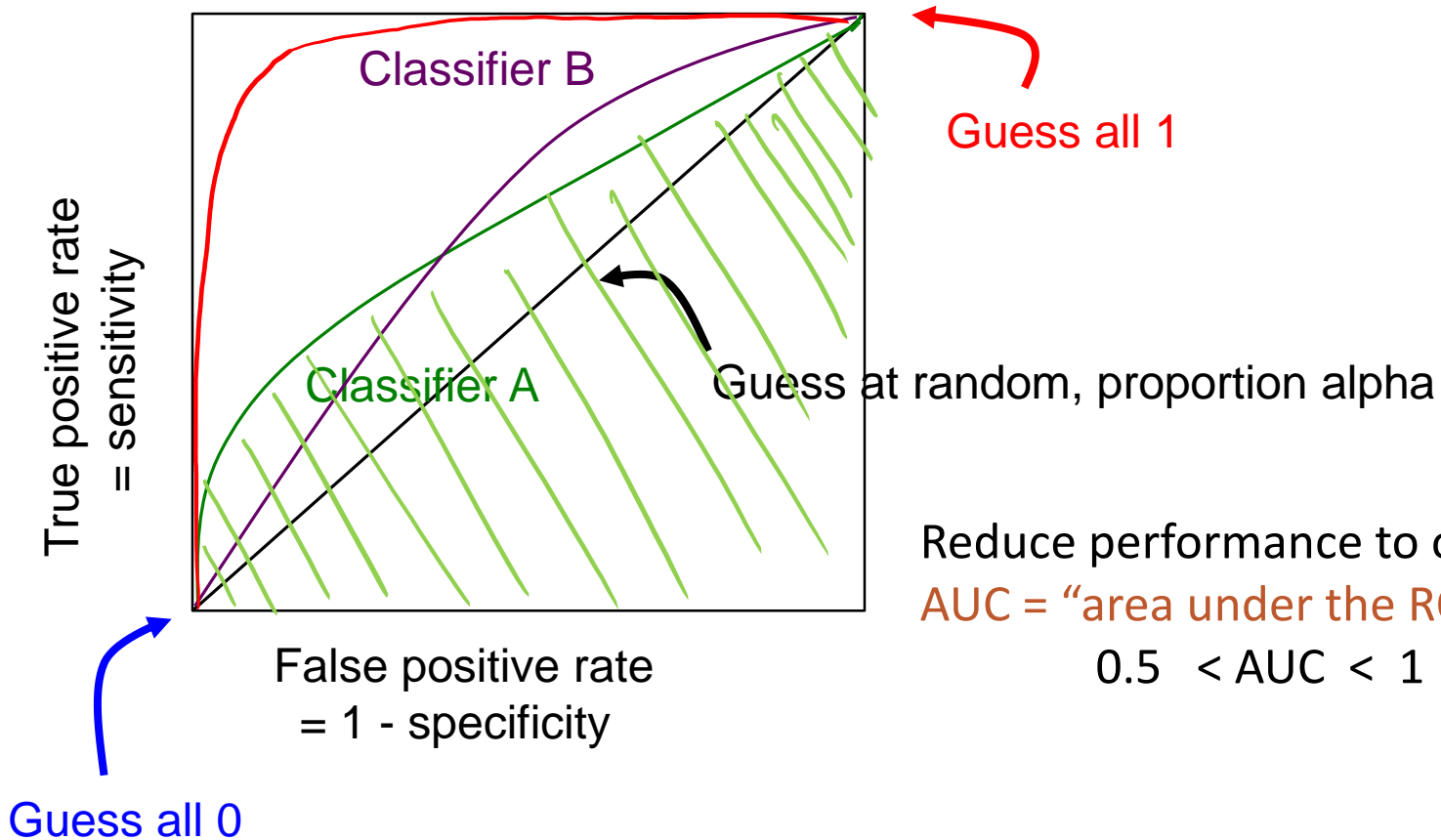*(expresses confidence in outcomes)*

Can use ROC curves for discriminative models also:
◦ Some notion of confidence, but doesn't correspond to a probability
◦ In our code: "predictSoft"  (vs. hard prediction, "predict")

```
>> learner = gaussianBayesClassify(X,Y);   % build a classifier
>> Ysoft = predictSoft(learner, X);        %  M x C matrix of confidences
>> plotSoftClassify2D(learner,X,Y);        %  shaded confidence plot
```

# ROC Curves

Characterize performance as we vary our confidence threshold?



True positive rate = sensitivity

False positive rate = 1 - specificity

Classifier B

Classifier A

Guess all 1

Guess at random, proportion alpha

Guess all 0

Reduce performance to one number?

AUC = "area under the ROC curve"

0.5  < AUC  < 1

# Machine Learning

Bayes Error Wrapup

**Gaussian Bayes Models**

Linear Regression: Definition and Cost

Gradient Descent Algorithms
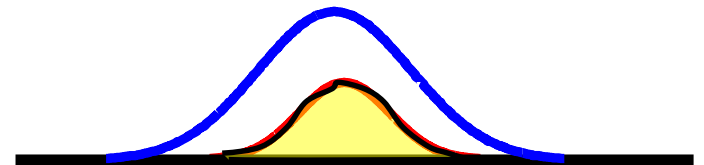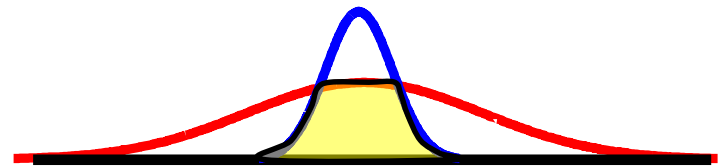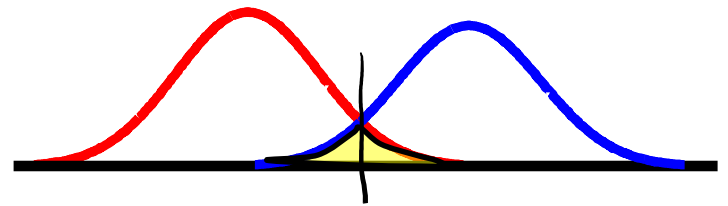
# Gaussian models

"Bayes optimal" decision
◦ Choose most likely class

Decision boundary
◦ Places where probabilities equal

What shape is the boundary?

# Gaussian models

Bayes optimal decision boundary
◦ p(y=0 | x) = p(y=1 | x)
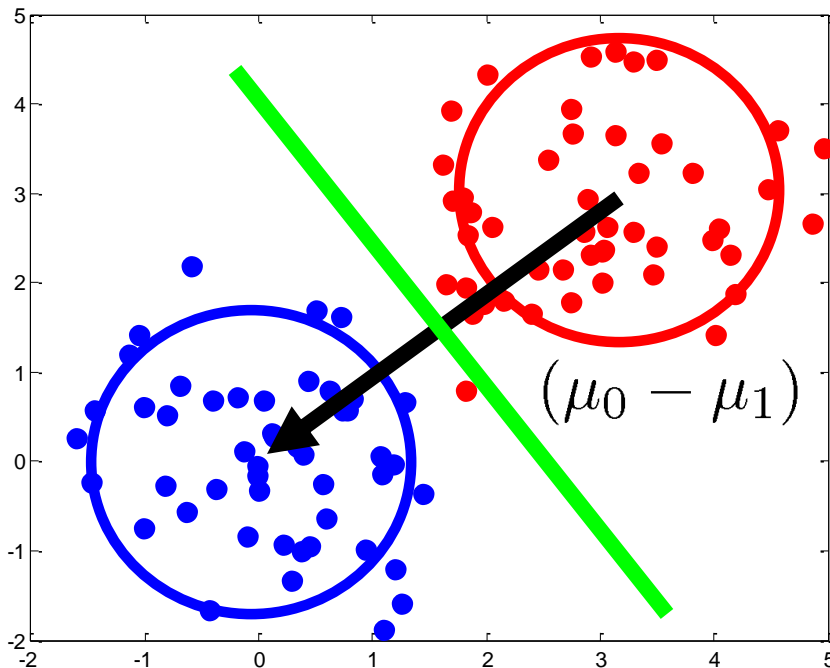◦ Transition point between p(y=0|x) >/< p(y=1|x)

$$p(x|y=1)$$

$$\mathcal{N}(\underline{x} \; ; \; \underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2}} |\Sigma|^{-1/2} \exp\left\{ -\frac{1}{2}(\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu}) \right\}$$

$$0 \begin{array}{c} < \\ > \end{array} \log \frac{p(x|y=0)}{p(x|y=1)} \frac{p(y=0)}{p(y=1)} = \log \frac{p(y=0)}{p(y=1)} - \frac{1}{2}\left( x\Sigma^{-1}x - 2\mu_0\Sigma^{-1}x + \mu_0^T\Sigma^{-1}\mu_0 \right)$$

$$+ \frac{1}{2}\left( x\Sigma^{-1}x - 2\mu_1\Sigma^{-1}x + \mu_1^T\Sigma^{-1}\mu_1 \right)$$

$$= (\mu_0 - \mu_1)^T \Sigma^{-1} x + \text{constants}$$

# Gaussian example

Spherical covariance: $\Sigma = \sigma^2 I$  $= (\mu_0 - \mu_1)^T \Sigma^{-1} x + constants$

Decision rule  $(\mu_0 - \mu_1)^T x \begin{array}{c} < \\ > \end{array} C$

$$C = .5(\mu_0^T \Sigma^{-1} \mu_0$$
$$- \mu_1^T \Sigma^{-1} \mu_1)$$
$$- \log \frac{p(y=0)}{p(y=1)}$$

$(\mu_0 - \mu_1)$

# Non-spherical Gaussian distributions

Equal covariances => still linear decision rule
- ◦ May be "modulated" by variance direction
- ◦ Scales; rotates (if correlated)

Example:
Variance
[ 3   0  ]
[ 0  .25 ]

# Class posterior probabilities
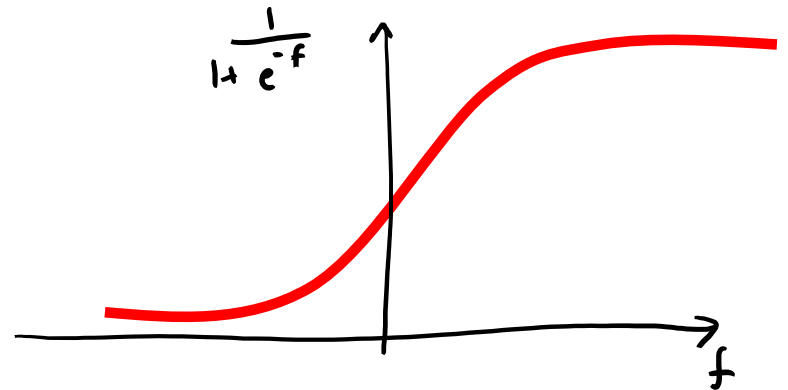
Consider comparing two classes
- p(x | y=0) * p(y=0)    vs    p(x | y=1) * p(y=1)
- Write probability of each class as
- p(y=0 | x) = p(y=0, x) / p(x)
-         = p(y=0, x) / ( p(y=0,x) + p(y=1,x) )   $= \dfrac{1}{1 + \dfrac{p(y=1,x)}{p(y=0,x)}}$
-     = 1 / (1  + exp( - $f$ ) )

- $f$ = log [ p(y=0, x) / p(y=1, x) ]

the logistic function, or logistic sigmoid

$\dfrac{1}{1 + e^{-f}}$

# Gaussian models

$$\mathcal{N}(\underline{x} \; ; \; \underline{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2}} |\Sigma|^{-1/2} \exp\left\{-\frac{1}{2}(\underline{x} - \underline{\mu})^T \Sigma^{-1}(\underline{x} - \underline{\mu})\right\}$$

$$0 \; \begin{matrix} < \\ > \end{matrix} \; \log \frac{p(x|y=0)}{p(x|y=1)} \frac{p(y=0)}{p(y=1)} = \; (\mu_0 - \mu_1)^T \Sigma^{-1} x + constants$$

$f$

Now we also know that the probability of each class is given by:
  p(y=0 | x) = Logistic( $f$ )  = Logistic(  $a^T$ x + b )

We'll see this form again soon…

# Machine Learning

Bayes Error Wrapup
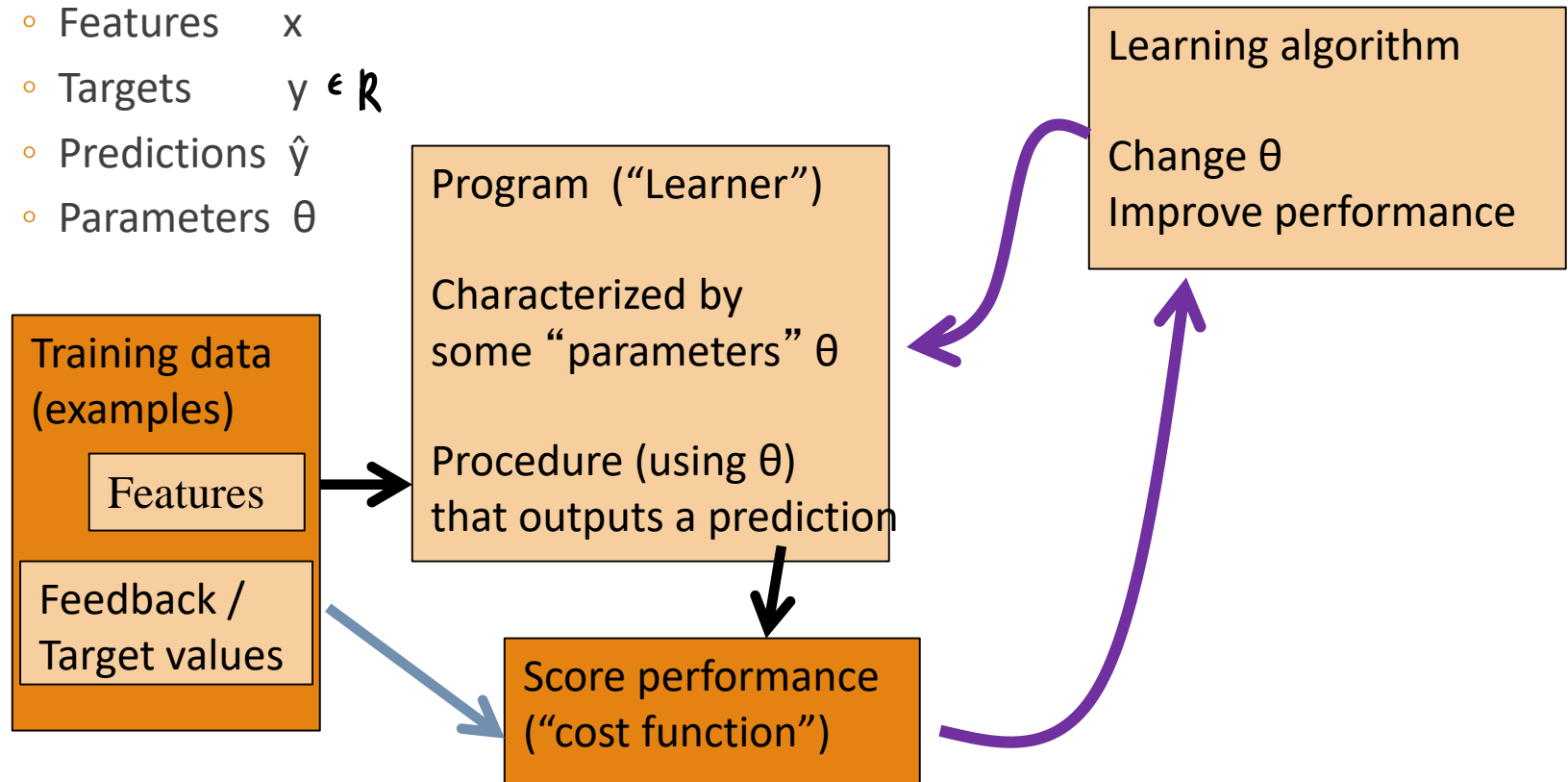
Gaussian Bayes Models

Linear Regression: Definition and Cost

Gradient Descent Algorithms

# Supervised learning

Notation

- ◦ Features     x
- ◦ Targets      y $\in \mathbb{R}$
- ◦ Predictions   ŷ
- ◦ Parameters   θ

**Program ("Learner")**

Characterized by
some "parameters" θ

Procedure (using θ)
that outputs a prediction

**Learning algorithm**

Change θ
Improve performance

**Training data
(examples)**

Features

Feedback /
Target values

**Score performance
("cost function")**

# Supervised learning

Notation
- Features      x
- Targets        y
- Predictions   ŷ
- Parameters   θ



Learning algorithm

Change θ
Improve performance

Program  ("Learner")

Characterized by
some "parameters" θ

Procedure (using θ)
that outputs a prediction

Training data
(examples)

Features

Feedback /
Target values

Score performance
("cost function")

# Linear regression



**"Predictor":**
Evaluate line:
$$r = \theta_0 + \theta_1 x_1$$

return r

Define form of function f(x) explicitly

Find a good f(x) within that family

# Notation

$$\hat{y}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_n x_n$$

Define feature $x_0 = 1$ (constant)
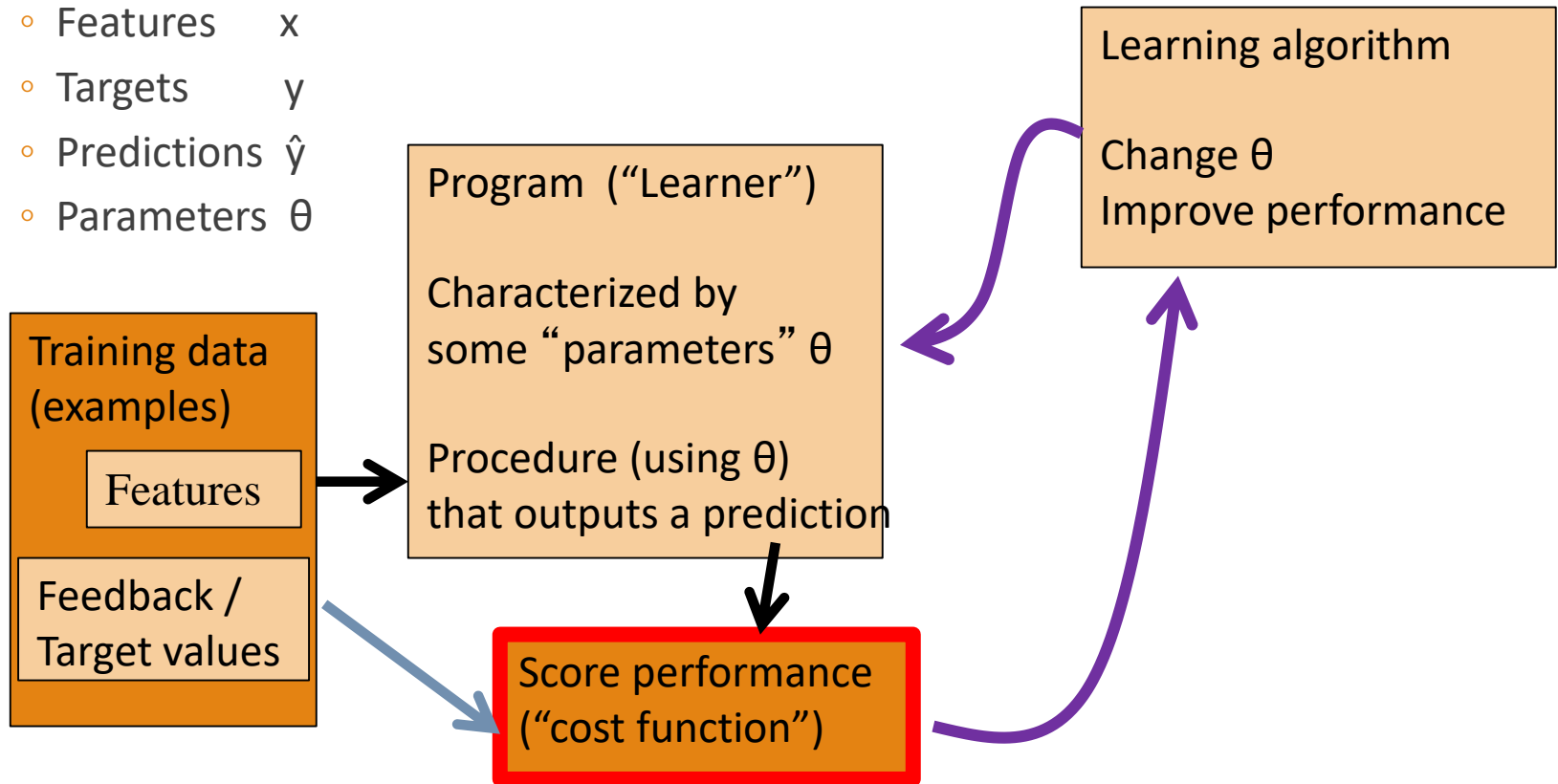
$$\hat{y}(x) = \theta x^T$$

$$\theta = [\theta_0, \theta_1, \ldots \theta_n] \quad 1 \times (n+1)$$

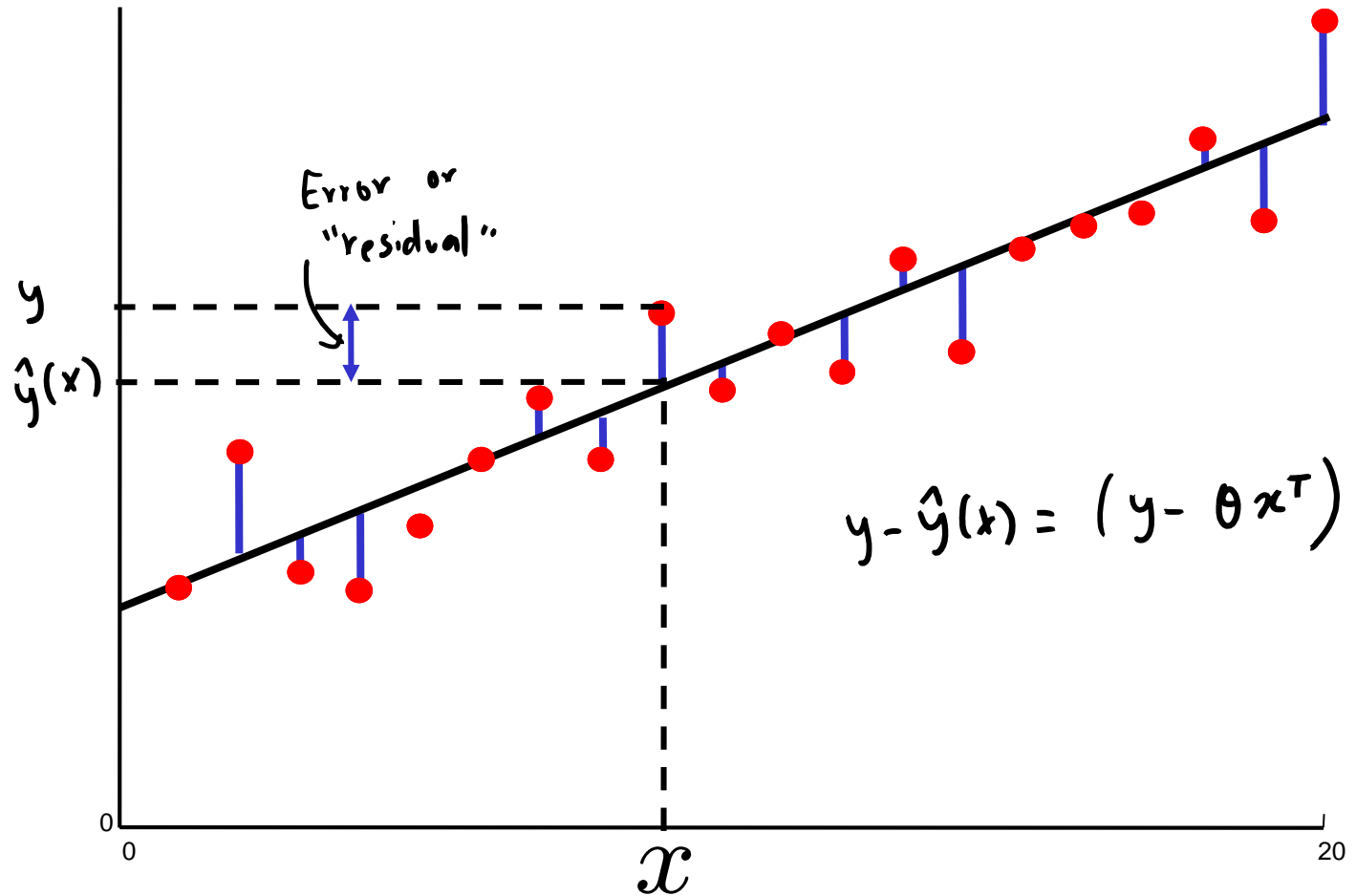$$x = [1, x_1, \ldots x_n] \quad 1 \times (n+1)$$

# Supervised learning

Notation
- Features      x
- Targets       y
- Predictions   ŷ
- Parameters   θ

**Training data (examples)**

Features

Feedback / Target values

**Program ("Learner")**

Characterized by some "parameters" θ

Procedure (using θ) that outputs a prediction

**Score performance ("cost function")**

**Learning algorithm**

Change θ
Improve performance

# Measuring error



Error or "residual"

$$y - \hat{y}(x) = \left( y - \theta x^{T} \right)$$

# Mean squared error

How can we quantify the error?

$$\text{MSE}, \quad J(\theta) = \frac{1}{m} \sum_j \left( y^{(j)} - \hat{y}(x^{(j)}) \right)^2$$

$$= \frac{1}{m} \sum_j \left( y^{(j)} - \theta \cdot x^{(j)T} \right)^2$$

Could choose something else, of course…

◦ Computationally convenient (more later)

◦ Measures the variance of the residuals

◦ Corresponds to likelihood under Gaussian model of "noise"

$$\mathcal{N}(y \; ; \; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{ -\frac{1}{2\sigma^2}(y - \mu)^2 \right\}$$

# MSE cost function

$$\text{MSE}, \quad J(\theta) = \frac{1}{m} \sum_{j} \left( y^{(j)} - \theta \cdot x^{(j)T} \right)^2$$

$$y = \left[ y^{(1)}, y^{(2)} \ldots y^{(m)} \right]^T$$

$$J(\theta) = \frac{1}{m} \left( y^T - \theta x^T \right) \left( y^T - \theta x^T \right)^T$$

$$X = \begin{bmatrix} x_0^{(1)} & x_1^{(1)} & \cdots & x_n^{(1)} \\ \vdots & \vdots & & \vdots \\ x_0^{(m)} & x_1^{(m)} & \cdots & x_n^{(m)} \end{bmatrix}$$

```python
# Python / NumPy:
e = Y - X.dot( theta.T );
J = e.T.dot( e ) / m   # = np.mean( e ** 2 )
```
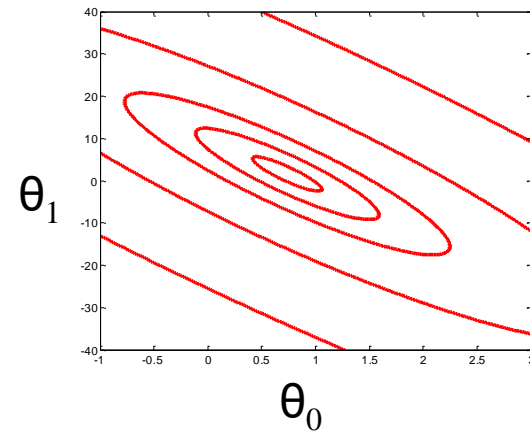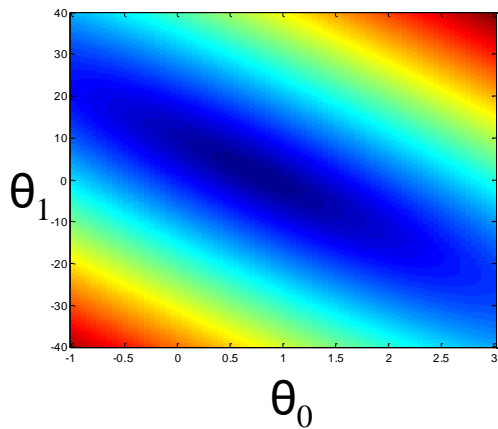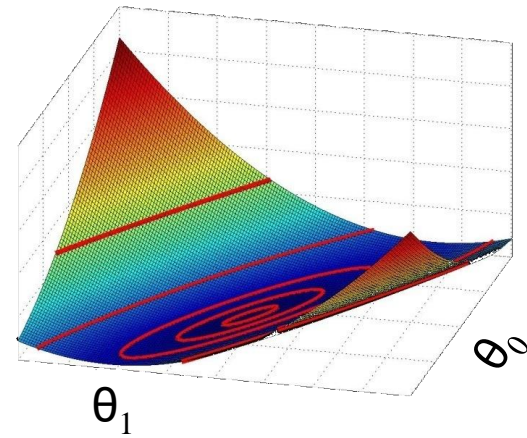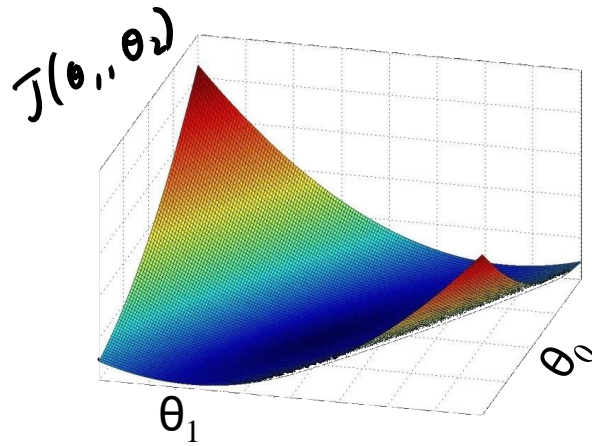
# Supervised learning

Notation

- Features    x
- Targets      y
- Predictions  ŷ
- Parameters  θ

Program  ("Learner")

Characterized by
some "parameters" θ

Procedure (using θ)
that outputs a prediction

Learning algorithm

Change θ
Improve performance

Training data
(examples)

Features

Feedback /
Target values

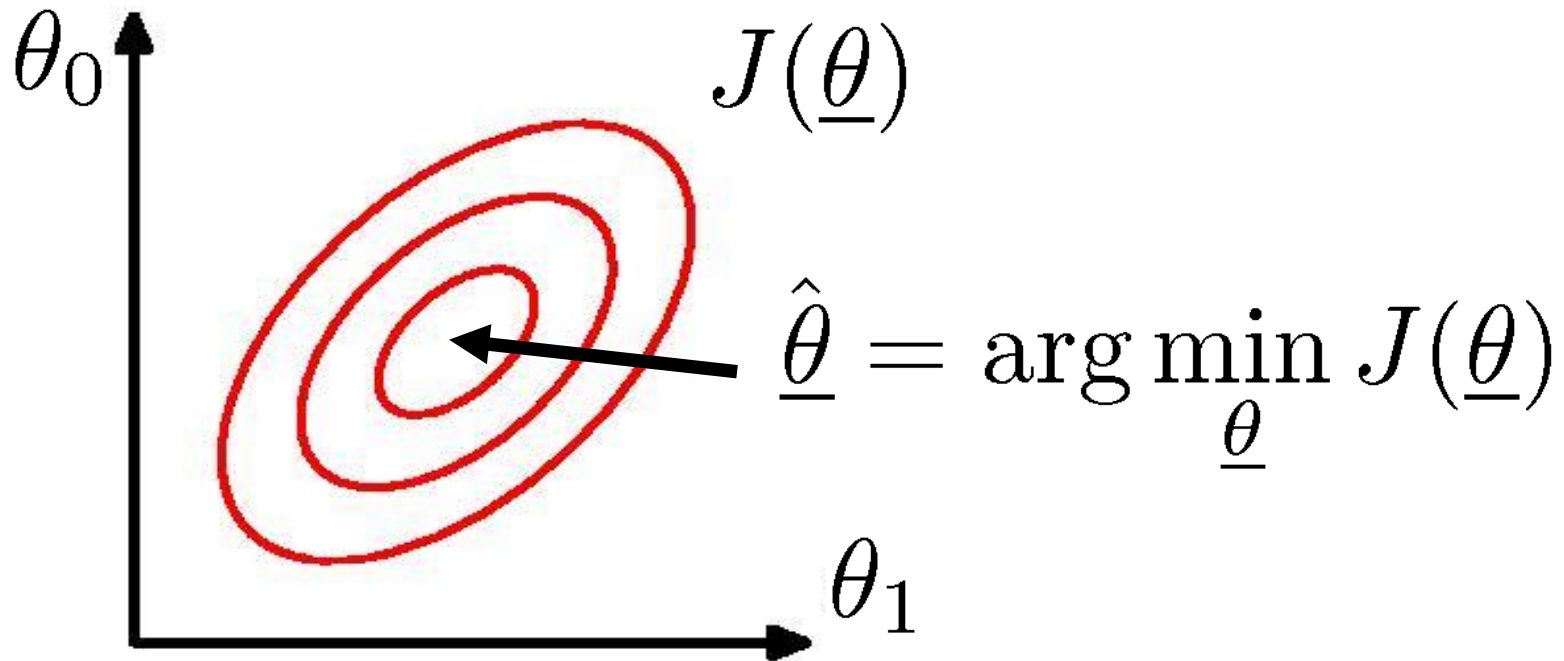Score performance
("cost function")

# Visualizing the cost function

# Finding good parameters

Want to find parameters which minimize our error…

Think of a cost "surface": error residual for that θ…

$$\theta_0$$

$$J(\underline{\theta})$$

$$\hat{\underline{\theta}} = \arg\min_{\underline{\theta}} J(\underline{\theta})$$

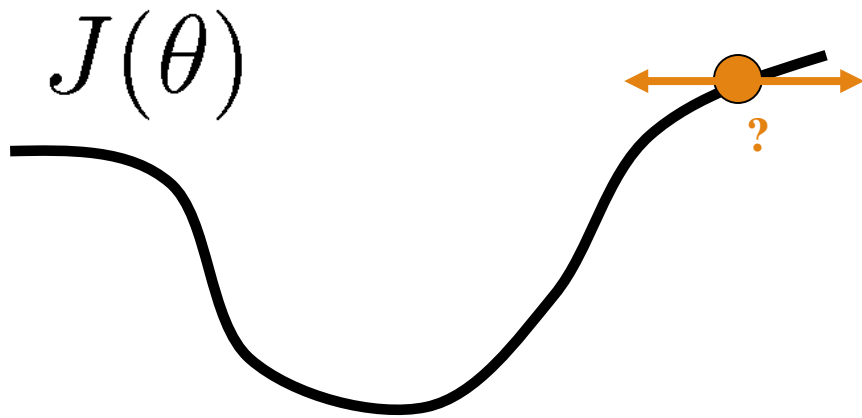$$\theta_1$$

# Machine Learning

Bayes Error Wrapup

Gaussian Bayes Models

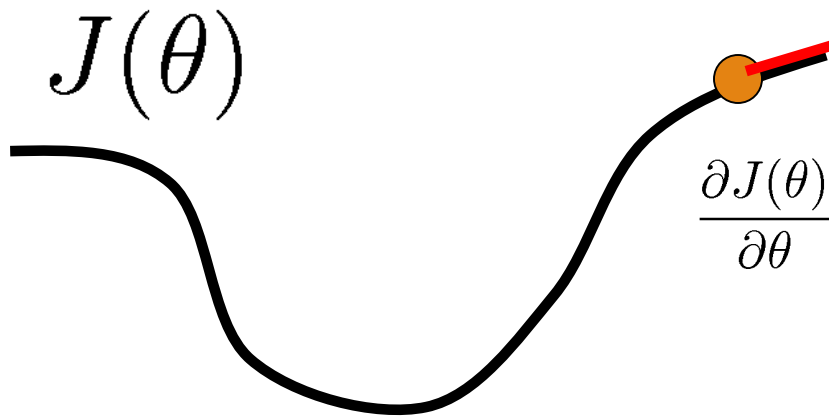Linear Regression: Definition and Cost

Gradient Descent Algorithms

# Gradient descent

$J(\theta)$

- How to change $\theta$ to improve $J(\theta)$?

- Choose a direction in which $J(\theta)$ is decreasing

# Gradient descent

$J(\theta)$

$\dfrac{\partial J(\theta)}{\partial \theta}$

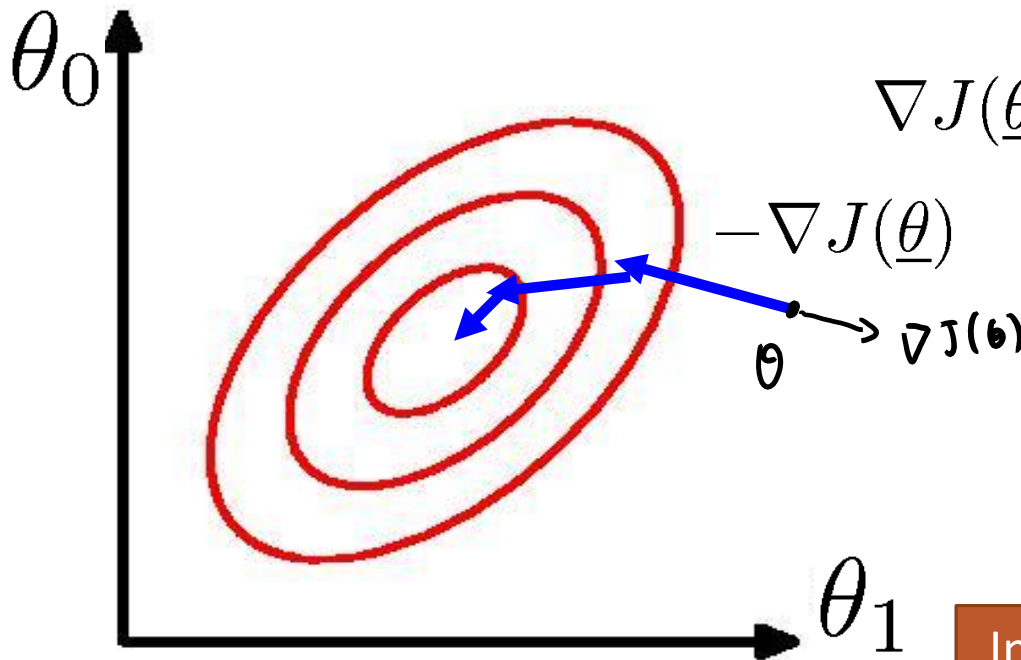- How to change θ to improve J(θ)?
- Choose a direction in which J(θ) is decreasing
- Derivative $\dfrac{\partial J(\theta)}{\partial \theta}$

- Positive => increasing
- Negative => decreasing

# Gradient descent in >2 dimensions



- Gradient vector

$$\nabla J(\underline{\theta}) = \left[ \frac{\partial J(\underline{\theta})}{\partial \theta_0} \quad \frac{\partial J(\underline{\theta})}{\partial \theta_1} \quad \cdots \right]$$

Indicates direction of steepest ascent
(negative = steepest descent)

# Gradient descent

Initialization

Step size
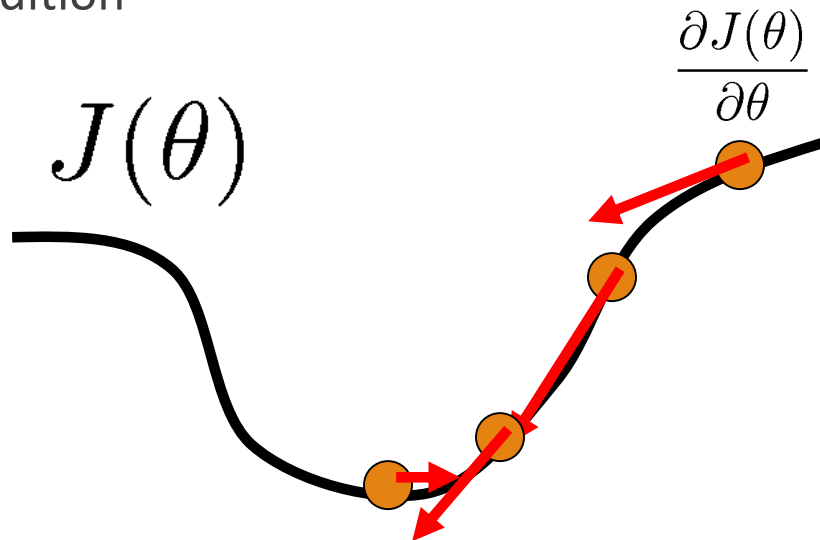◦ Can change as a function of iteration

Gradient direction

Stopping condition

Initialize θ — randomly (small)
                — $0_s$

Do{

  θ ← θ - α∇$_θ$J(θ)

} while (α‖∇$_θ$J‖ > ε)

$$\frac{\partial J(\theta)}{\partial \theta}$$

$$J(\theta)$$

# Gradient for the MSE

$$J(\theta) = \frac{1}{m} \sum_j \left( y^{(j)} - \theta_0 x_0^{(j)} - \theta_1 x_1^{(j)} \cdots - \theta_n x_n^{(j)} \right)^2$$

$$\underbrace{\phantom{y^{(j)} - \theta_0 x_0^{(j)} - \theta_1 x_1^{(j)} \cdots - \theta_n x_n^{(j)}}}_{\substack{e_j(\theta) \\ \text{(error on } j\text{)}}}$$

$$\nabla J(\theta) = \left[ \frac{\partial J(\theta)}{\partial \theta_0} \quad \frac{\partial J(\theta)}{\partial \theta_1} \quad \cdots \quad \frac{\partial J(\theta)}{\partial \theta_n} \right]$$

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{m} \frac{\partial}{\partial \theta_0} \sum_j \left( e_j(\theta) \right)^2$$

$$= \frac{1}{m} \sum_j 2 e_j(\theta) \frac{\partial e_j(\theta)}{\partial \theta_0}$$

$$\frac{\partial e_j(\theta)}{\partial \theta_0} = \frac{\partial y^{(j)}}{\partial \theta_0} - \frac{\partial \theta_0 x_0^{(j)}}{\partial \theta_0} - \frac{\partial \theta_1 x_1^{(j)}}{\partial \theta_0} + \cdots$$

$$= - x_0^{(j)}$$

# Upcoming…

**Misc.**
- Lot of activity on Piazza
- You have been added to Gradescope

**Homework**
- Homework 1 due tonight
- Homework 2 released tonight
- HW2 Due: October 19, 2017