



# Rensselaer

## Lecture 20: Regression

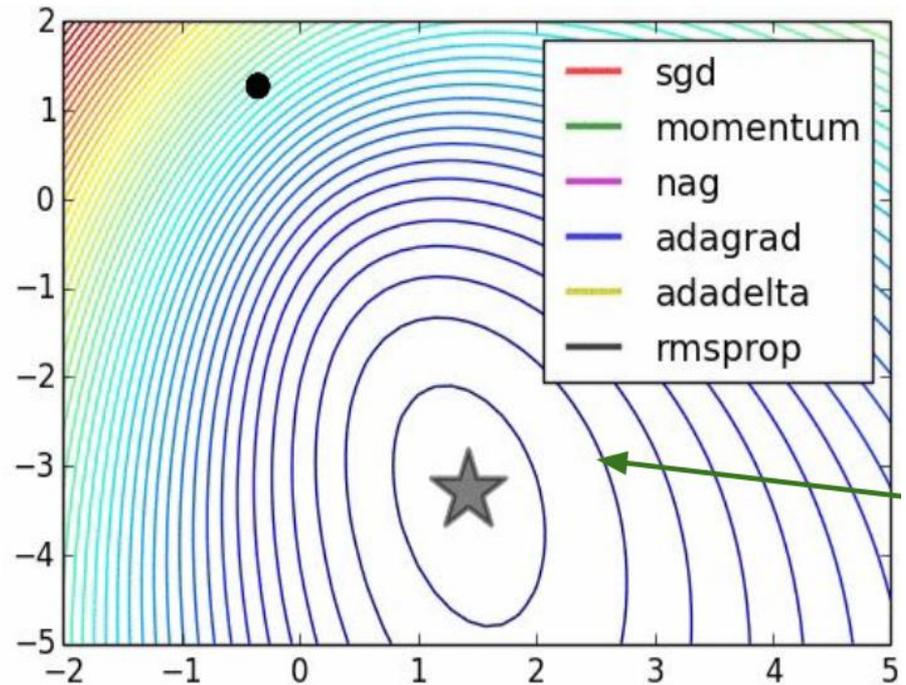
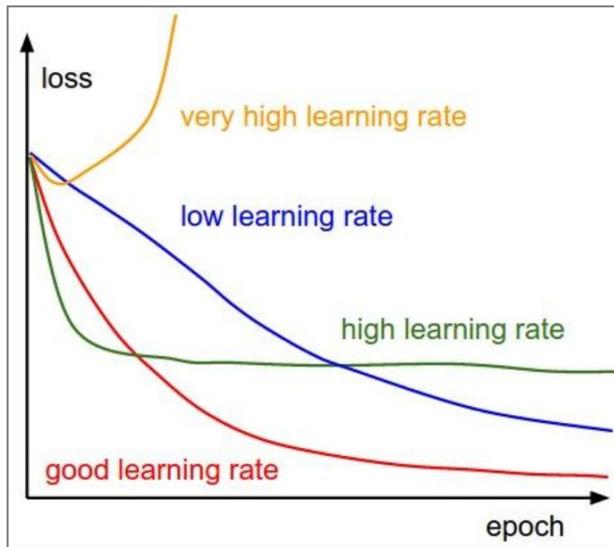
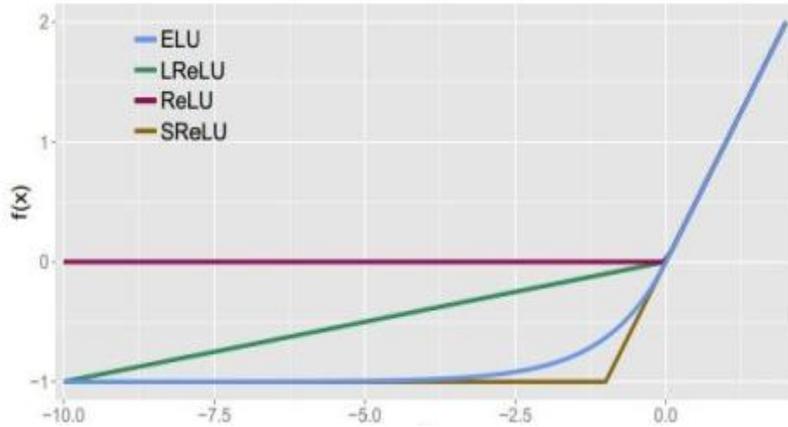
Dr. Chengjiang Long

Computer Vision Researcher at Kitware Inc.

Adjunct Professor at RPI.

Email: [longc3@rpi.edu](mailto:longc3@rpi.edu)

# Recap Previous Lecture



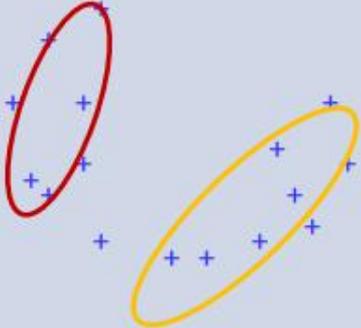
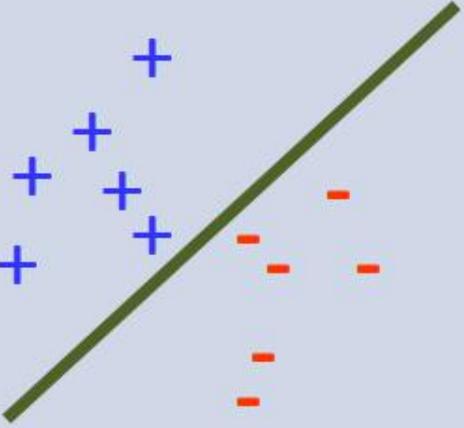
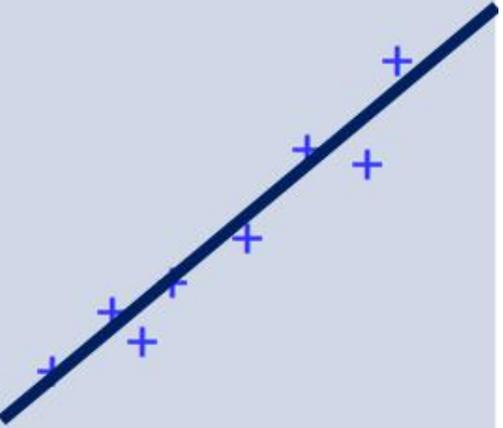
# Outline

- Regression Overview
- Linear Regression
- Support Vector Regression
- Logistic Regression
- Deep Neural Network for Regression

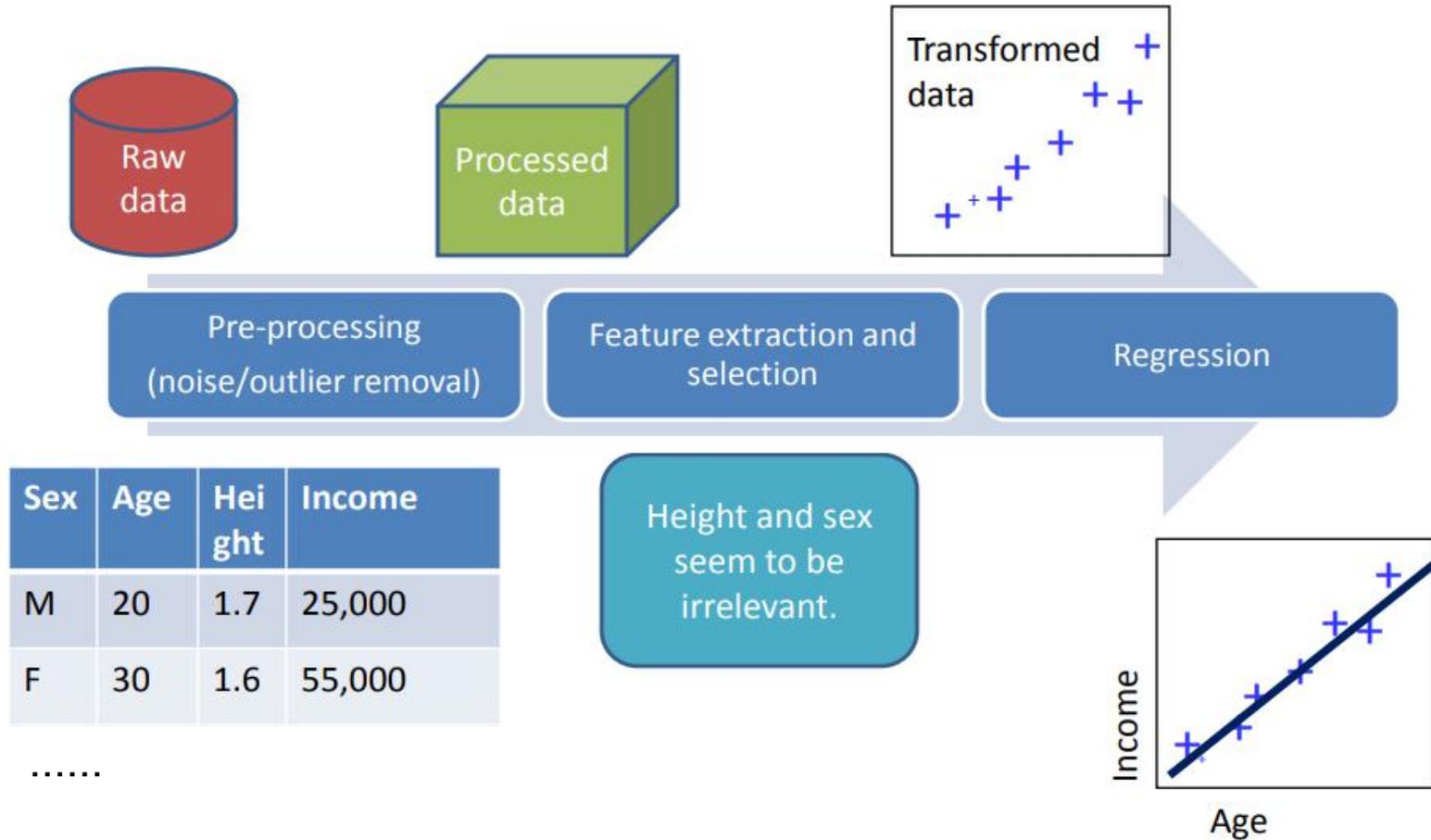
# Outline

- **Regression Overview**
- Linear Regression
- Support Vector Regression
- Logistic Regression
- Deep Neural Network for Regression

# Regression Overview

CLUSTERING	CLASSIFICATION	REGRESSION (THIS TALK)
		
<p>K-means</p> <p>Hierarchy clustering</p> <p>Gaussian Mixture Model</p>	<ul style="list-style-type: none"> <li>• Decision tree</li> <li>• Linear Discriminant Analysis</li> <li>• Neural Networks</li> <li>• Support Vector Machines</li> <li>• Boosting</li> </ul>	<ul style="list-style-type: none"> <li>• Linear Regression</li> <li>• Support Vector Regression</li> <li>• Logistic Regression</li> <li>• Neural Networks</li> </ul>
<p>Group data based on their characteristics</p>	<p>Separate data based on their labels</p>	<p>Find a model that can explain the output given the input</p>

# One Example



# Evaluation Metrics

- Root mean-square error (RMSE)
  - RMSE is a popular formula to measure the error rate of a regression model. However, it can only be compared between models whose errors are measured in the same units.

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (p_i - a_i)^2}{n}}$$

$a$  = actual target

$p$  = predicted target

- Mean absolute error (MAE)
  - MAE has the same unit as the original data, and it can only be compared between models whose errors are measured in the same units. It is usually similar in magnitude to RMSE, but slightly smaller.

$$MAE = \frac{\sum_{i=1}^n |p_i - a_i|}{n}$$

# Evaluation Metrics

- Relative Squared Error (RSE)

- Unlike RMSE, the relative squared error (RSE) can be compared between models whose errors are measured in the different units.

$$RSE = \frac{\sum_{i=1}^n (p_i - a_i)^2}{\sum_{i=1}^n (\bar{a} - a_i)^2}$$

mean value

- Relative Absolute Error (RAE)

- Like RSE , the relative absolute error (RAE) can be compared between models whose errors are measured in the different units..

$$RAE = \frac{\sum_{i=1}^n |p_i - a_i|}{\sum_{i=1}^n |\bar{a} - a_i|}$$

# Outline

- Regression Overview
- **Linear Regression**
- Support Vector Regression
- Logistic Regression
- Deep Neural Network for Regression

# Linear Regression

- Given data with  $n$  dimensional variables and 1 target-variable (real number)

$$\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$$

where  $\mathbf{x} \in \mathcal{R}^n, y \in \mathcal{R}$

- The objective: Find a function  $f$  that returns the best fit.

$$f : \mathcal{R}^n \rightarrow \mathcal{R}$$

- Assume that the relationship between  $X$  and  $y$  is approximately linear. The model can be represented as ( $w$  represents coefficients and  $b$  is an intercept)

$$f(w_1, \dots, w_n, b) = y = \mathbf{w} \cdot \mathbf{x} + b + \varepsilon$$

# Linear Regression

- To find the best fit, we minimize the sum of squared errors → Least square estimation

$$\min \sum_{i=1}^m (y_i - \hat{y}_i)^2 = \sum_{i=1}^m (y_i - (\mathbf{w} \cdot \mathbf{x}_i + b))^2$$

- The solution can be found by solving

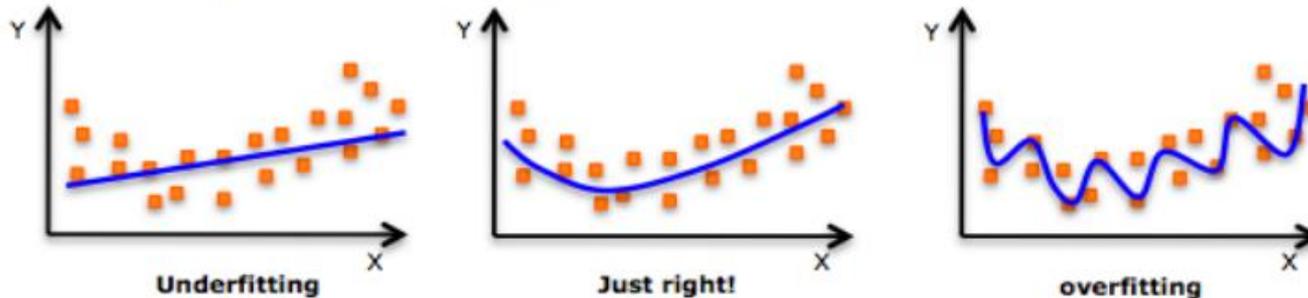
$$\hat{\mathbf{w}} = (X^T X)^{-1} X^T Y$$

(By taking the derivative of the above objective function w.r.t. )

- In MATLAB, the back-slash operator computes a least square solution.

# Linear Regression

$$\min \sum_{i=1}^m (y_i - \hat{y}_i)^2 = \sum_{i=1}^m (y_i - (\hat{\mathbf{w}} \cdot \mathbf{x}_i + \hat{b}))^2$$



To avoid over-fitting, a regularization term can be introduced (minimize a magnitude of  $w$ )

– LASSO: 
$$\min \sum_{i=1}^m (y_i - \mathbf{w} \cdot \mathbf{x}_i - b)^2 + C \sum_{j=1}^n |w_j|$$

– Ridge regression: 
$$\min \sum_{i=1}^m (y_i - \mathbf{w} \cdot \mathbf{x}_i - b)^2 + C \sum_{j=1}^n |\mathbf{w}_j^2|$$

# Outline

- Regression Overview
- Linear Regression
- **Support Vector Regression**
- Logistic Regression
- Deep Neural Network for Regression

# Support Vector Regression

- Find a function,  $f(x)$ , with at most  $\varepsilon$  -deviation from the target  $y$

The problem can be written as a convex optimization problem

$$\min \frac{1}{2} \|\mathbf{w}\|^2$$

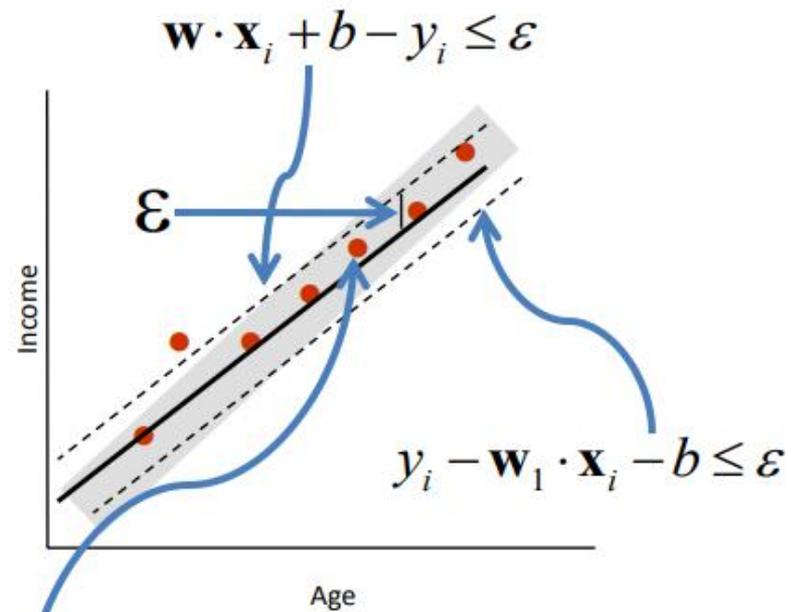
$$s.t. \ y_i - \mathbf{w}_1 \cdot \mathbf{x}_i - b \leq \varepsilon;$$

$$\mathbf{w}_1 \cdot \mathbf{x}_i + b - y_i \leq \varepsilon;$$

C: trade off the complexity

What if the problem is not feasible?

We can introduce slack variables  
(similar to soft margin loss function).

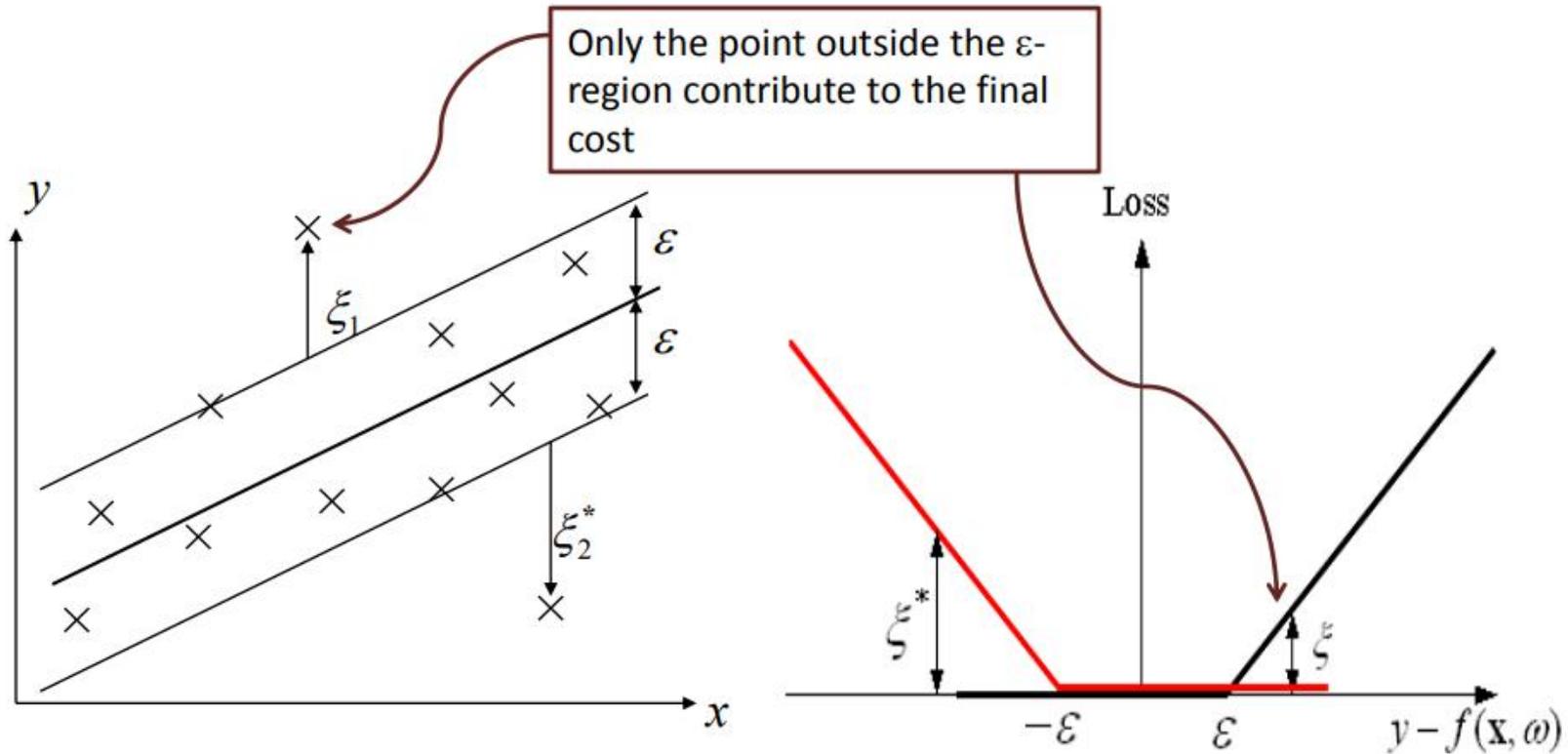


We do not care about errors as long as they are less than  $\varepsilon$

# Support Vector Regression

Assume linear parameterization

$$f(\mathbf{x}, \omega) = \mathbf{w} \cdot \mathbf{x} + b$$



$$L_\varepsilon(y, f(\mathbf{x}, \omega)) = \max(|y - f(\mathbf{x}, \omega)| - \varepsilon, 0)$$

# Soft margin

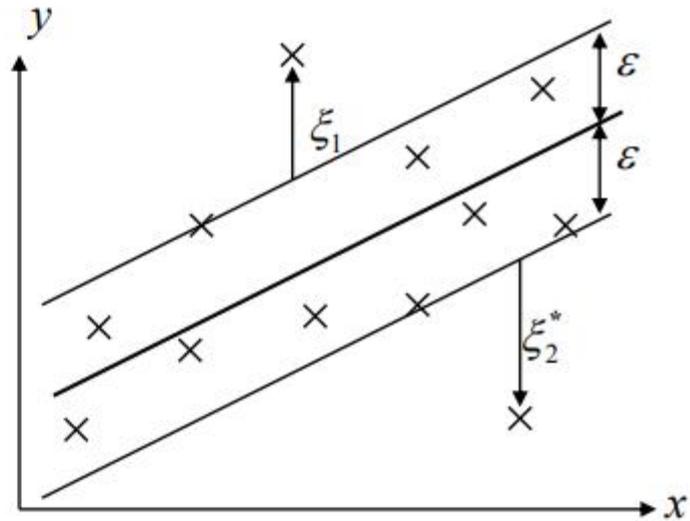
Given training data

$$(\mathbf{x}_i, y_i) \quad i = 1, \dots, m$$

Minimize

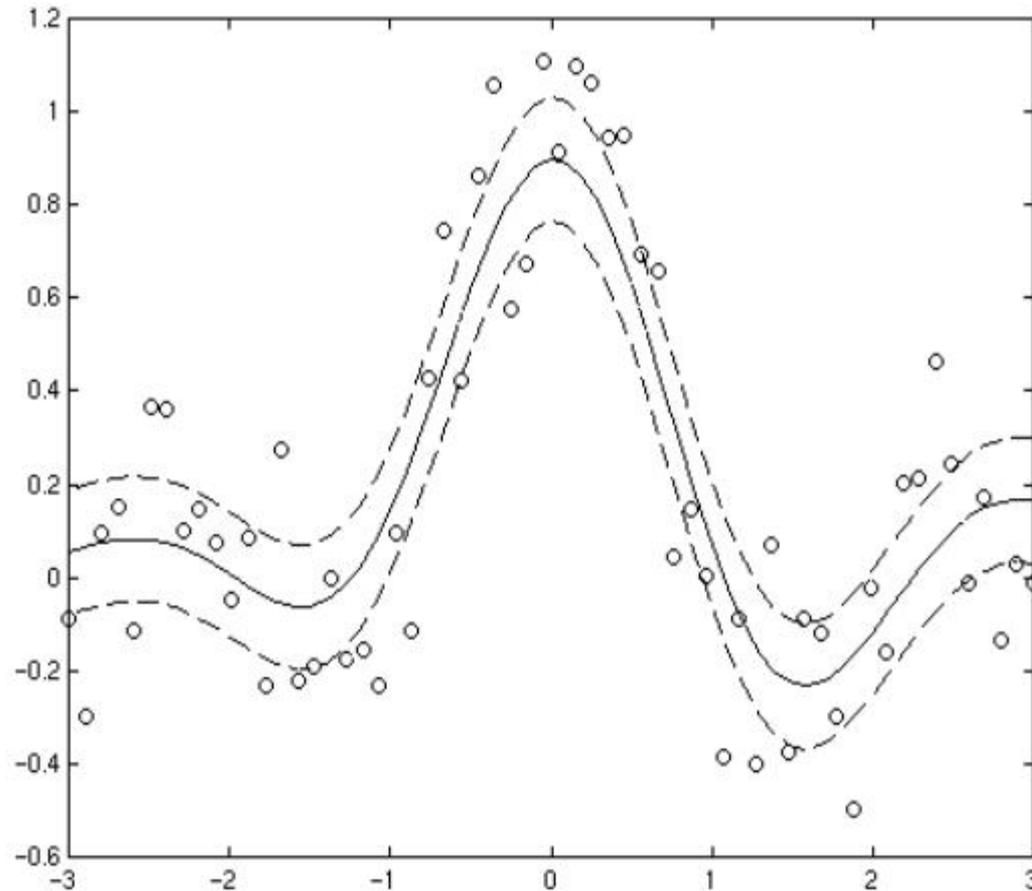
$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*)$$

Under constraints



$$\begin{cases} y_i - (\mathbf{w} \cdot \mathbf{x}_i) - b \leq \epsilon + \xi_i \\ (\mathbf{w} \cdot \mathbf{x}_i) + b - y_i \leq \epsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0, i = 1, \dots, m \end{cases}$$

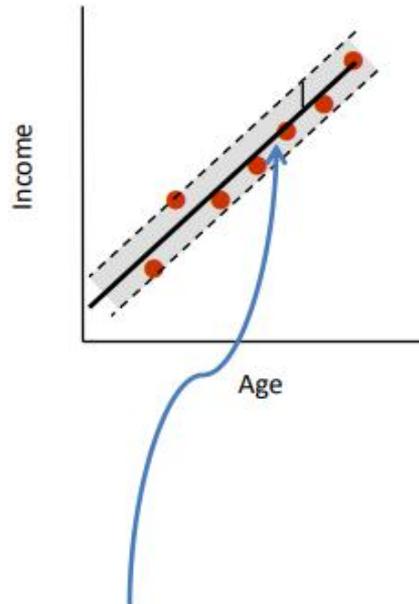
# How about a non-linear case?



# Linear versus Non-linear SVR

- **Linear case**

$$f : \text{age} \rightarrow \text{income}$$

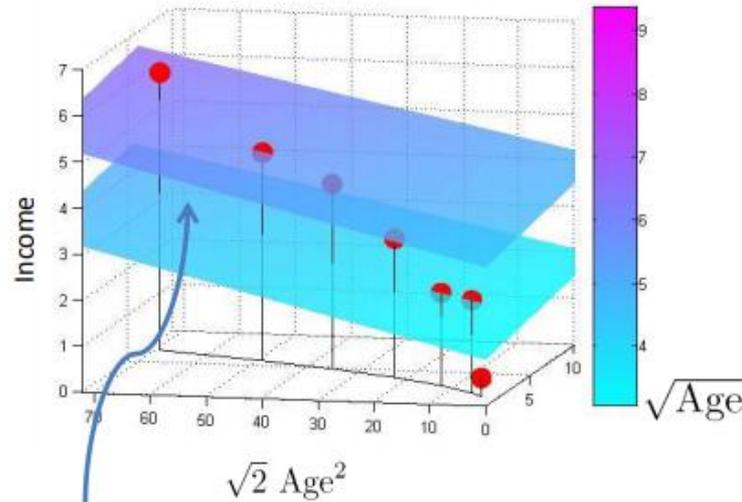


$$y_i = \mathbf{w}_1 \cdot \mathbf{x}_i + b$$

- **Non-linear case**

- Map data into a higher dimensional space, e.g.,

$$f : (\sqrt{\text{age}}, \sqrt{2\text{age}^2}) \rightarrow \text{income}$$



$$y_i = \mathbf{w}_1 \sqrt{\mathbf{x}_i} + \mathbf{w}_2 \sqrt{2\mathbf{x}_i^2} + b$$

# Dual problem

- Primal

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*)$$

$$s.t. \begin{cases} y_i - (\mathbf{w} \cdot \mathbf{x}_i) - b \leq \varepsilon + \xi_i \\ (\mathbf{w} \cdot \mathbf{x}_i) + b - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0, i = 1, \dots, m \end{cases}$$

- Dual

$$\max \begin{cases} \frac{1}{2} \sum_{i,j=1}^m (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \langle x_i, x_j \rangle \\ -\varepsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) + \sum_{i=1}^m y_i (\alpha_i - \alpha_i^*) \end{cases}$$

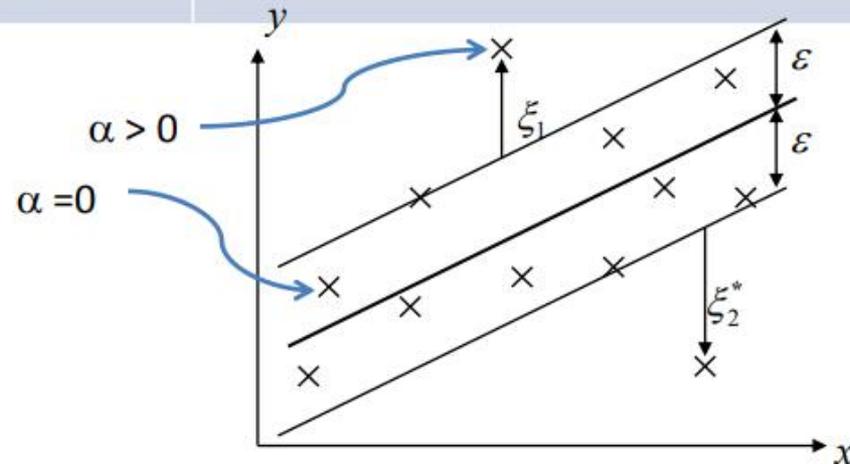
$$s.t. \sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0; 0 \leq \alpha_i, \alpha_i^* \leq C$$

Primal variables:  $w$  for each feature dim

Dual variables:  $\alpha, \alpha^*$  for each data point

Complexity: the dim of the input space

Complexity: Number of support vectors



# Kernel trick

- Linear:  $\langle x, y \rangle$
- Non-linear:  $\langle \varphi(x), \varphi(y) \rangle = K(x, y)$

Note: No need to compute the mapping function,  $\varphi(\cdot)$ , explicitly. Instead, we use the kernel function.

## Commonly used kernels:

- Polynomial kernels:  $K(x, y) = (x^T y + 1)^d$

- Radial basis function (RBF) kernels:

$$K(x, y) = \exp\left(-\frac{1}{2\sigma^2} \|x - y\|^2\right)$$

Note: for RBF kernel,  $\dim(\varphi(\cdot))$  is infinite

# Dual problem for non-linear case

- Primal

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m (\xi_i + \xi_i^*)$$

$$s.t. \begin{cases} y_i - (\mathbf{w} \cdot \varphi(\mathbf{x}_i)) - b \leq \varepsilon + \xi_i \\ (\mathbf{w} \cdot \varphi(\mathbf{x}_i)) + b - y_i \leq \varepsilon + \xi_i^* \\ \xi_i, \xi_i^* \geq 0, i = 1, \dots, m \end{cases}$$

- Dual

$$\max \begin{cases} \frac{1}{2} \sum_{i,j=1}^m (\alpha_i - \alpha_i^*)(\alpha_j - \alpha_j^*) \langle \varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j) \rangle \\ -\varepsilon \sum_{i=1}^m (\alpha_i + \alpha_i^*) + \sum_{i=1}^m y_i (\alpha_i - \alpha_i^*) \end{cases}$$

$K(\mathbf{x}_i, \mathbf{x}_j)$

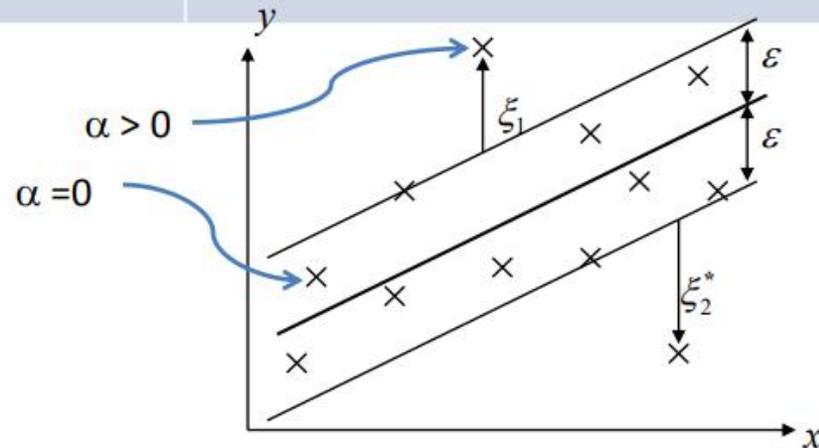
$$s.t. \sum_{i=1}^m (\alpha_i - \alpha_i^*) = 0; 0 \leq \alpha_i, \alpha_i^* \leq C$$

Primal variables:  $\mathbf{w}$  for each feature dim

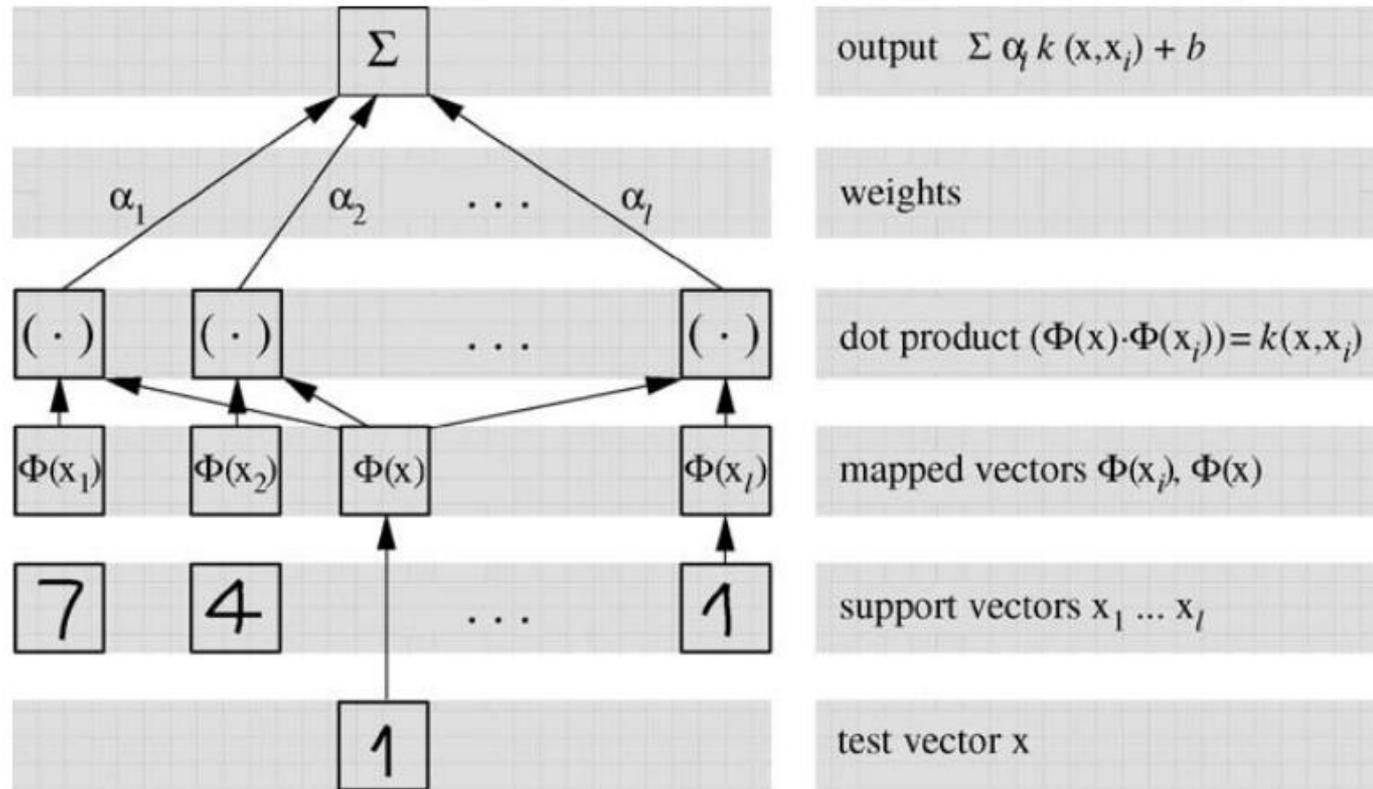
Dual variables:  $\alpha, \alpha^*$  for each data point

Complexity: the dim of the input space

Complexity: Number of support vectors



# Architecture of a regression machine



[Alex J. Smola et al. A tutorial on support vector regression, 2004.]  
URL: <https://alex.smola.org/papers/2004/SmoSch04.pdf>

# Outline

- Regression Overview
- Linear Regression
- Support Vector Regression
- **Logistic Regression**
- Deep Neural Network for Regression

# Logistic Regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)

$h_{\theta}(\mathbf{x})$  should give  $p(y = 1 \mid \mathbf{x}; \theta)$

– Want  $0 \leq h_{\theta}(\mathbf{x}) \leq 1$

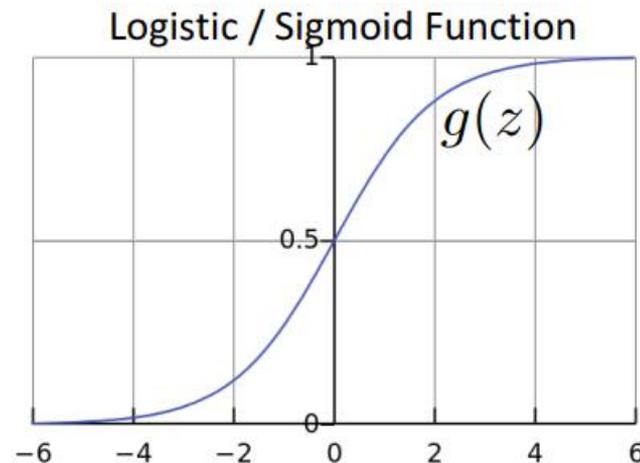
Can't just use linear regression with a threshold

- Logistic regression model:

$$h_{\theta}(\mathbf{x}) = g(\theta^T \mathbf{x})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$



# Interpretation of Hypothesis Output

$$h_{\theta}(\mathbf{x}) = \text{estimated } p(y = 1 \mid \mathbf{x}; \theta)$$

Example: Cancer diagnosis from tumor size

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$

$$h_{\theta}(\mathbf{x}) = 0.7$$

→ Tell patient that 70% chance of tumor being malignant

Note that:  $p(y = 0 \mid \mathbf{x}; \theta) + p(y = 1 \mid \mathbf{x}; \theta) = 1$

Therefore,  $p(y = 0 \mid \mathbf{x}; \theta) = 1 - p(y = 1 \mid \mathbf{x}; \theta)$

# Another Interpretation

- Equivalently, logistic regression assumes that

$$\log \frac{p(y = 1 \mid \mathbf{x}; \boldsymbol{\theta})}{p(y = 0 \mid \mathbf{x}; \boldsymbol{\theta})} = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

odds of  $y = 1$

**Side Note:** the odds in favor of an event is the quantity  $p / (1 - p)$ , where  $p$  is the probability of the event

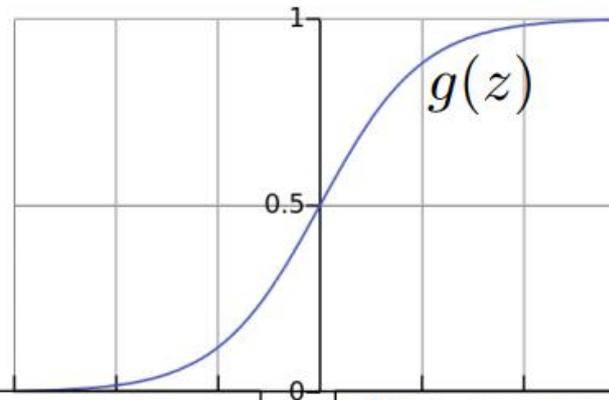
E.g., If I toss a fair dice, what are the odds that I will have a 6?

- In other words, logistic regression assumes that the log odds is a linear function of  $\mathbf{x}$

# Logistic Regression

$$h_{\theta}(\mathbf{x}) = g(\theta^{\top} \mathbf{x})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

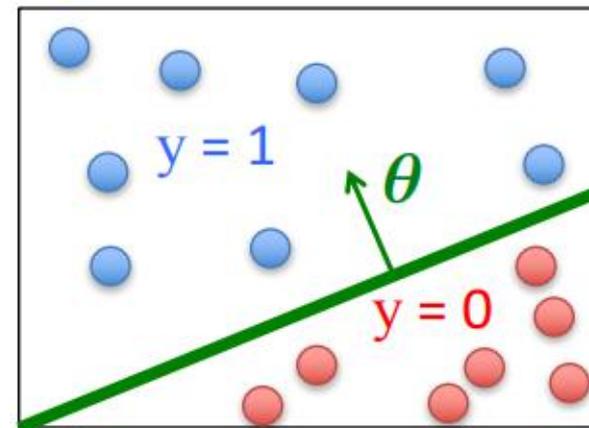


$\theta^{\top} \mathbf{x}$  should be large negative values for negative instances

$\theta^{\top} \mathbf{x}$  should be large positive values for positive instances

- Assume a threshold and...

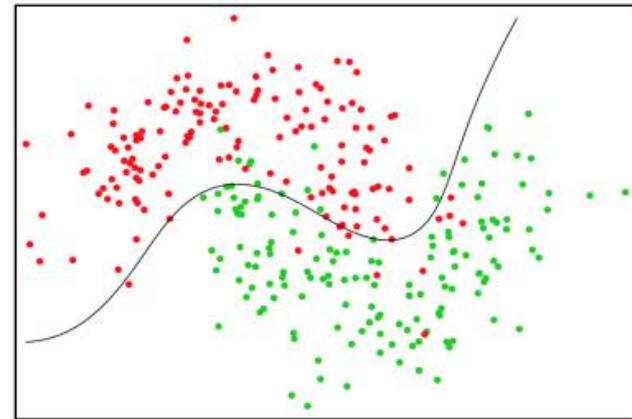
- Predict  $y = 1$  if  $h_{\theta}(\mathbf{x}) \geq 0.5$
- Predict  $y = 0$  if  $h_{\theta}(\mathbf{x}) < 0.5$



# Non-Linear Decision Boundary

- Can apply basis function expansion to features, same as with linear regression?

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ \vdots \end{bmatrix}$$



# Logistic Regression

Given  $\left\{ \left( \mathbf{x}^{(1)}, y^{(1)} \right), \left( \mathbf{x}^{(2)}, y^{(2)} \right), \dots, \left( \mathbf{x}^{(n)}, y^{(n)} \right) \right\}$

where  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ ,  $y^{(i)} \in \{0, 1\}$

Model:  $h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^T \mathbf{x})$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^T = \begin{bmatrix} 1 & x_1 & \dots & x_d \end{bmatrix}$$

# Logistic Regression Objective Function

- Can't just use squared loss as in linear regression

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

- Using the logistic regression model

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

results in a non-convex optimization

# Deriving the Cost Function via Maximum Likelihood Estimation

- Likelihood of data is given by:  $l(\boldsymbol{\theta}) = \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$

- So, looking for the  $\boldsymbol{\theta}$  that maximizes the likelihood

$$\boldsymbol{\theta}_{\text{MLE}} = \arg \max_{\boldsymbol{\theta}} l(\boldsymbol{\theta}) = \arg \max_{\boldsymbol{\theta}} \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta})$$

- Can take the log without changing the solution:

$$\begin{aligned} \boldsymbol{\theta}_{\text{MLE}} &= \arg \max_{\boldsymbol{\theta}} \log \prod_{i=1}^n p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^n \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \end{aligned}$$

# Deriving the Cost Function via Maximum Likelihood Estimation

- Expand as follows:

$$\begin{aligned}\theta_{\text{MLE}} &= \arg \max_{\theta} \sum_{i=1}^n \log p(y^{(i)} | \mathbf{x}^{(i)}; \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^n \left[ y^{(i)} \log p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \theta) + (1 - y^{(i)}) \log (1 - p(y^{(i)} = 1 | \mathbf{x}^{(i)}; \theta)) \right]\end{aligned}$$

- Substitute in model, and take negative to yield

**Logistic regression objective:**

$$\begin{aligned}\min_{\theta} J(\theta) \\ J(\theta) &= - \sum_{i=1}^n \left[ y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(\mathbf{x}^{(i)})) \right]\end{aligned}$$

# Intuition Behind the Objective

$$J(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$

- Cost of a single instance:

$$\text{cost}(h_{\boldsymbol{\theta}}(\mathbf{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

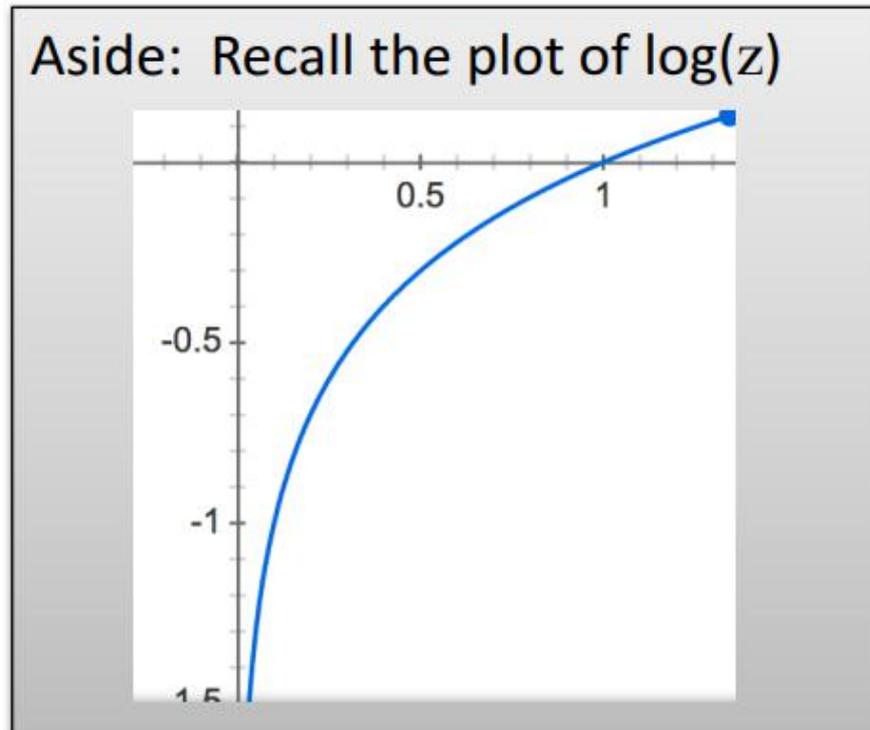
- Can re-write objective function as

$$J(\boldsymbol{\theta}) = \sum_{i=1}^n \text{cost}(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)})$$

Compare to linear regression:  $J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$

# Intuition Behind the Objective

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

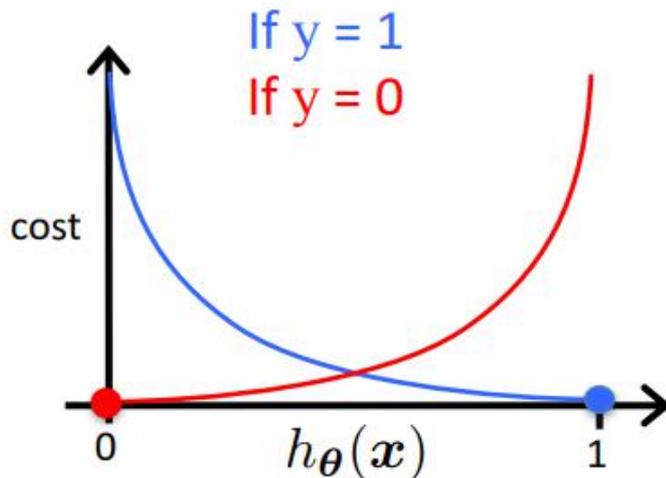


# Intuition Behind the Objective

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

If  $y = 0$

- Cost = 0 if prediction is correct
- As  $(1 - h_{\theta}(\mathbf{x})) \rightarrow 0$ ,  $\text{cost} \rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties



# Regularized Logistic Regression

$$J(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$

- We can regularize logistic regression exactly as before

$$\begin{aligned} J_{\text{regularized}}(\boldsymbol{\theta}) &= J(\boldsymbol{\theta}) + \lambda \sum_{j=1}^d \theta_j^2 \\ &= J(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}_{[1:d]}\|_2^2 \end{aligned}$$

# Gradient Descent for Logistic Regression

$$J_{\text{reg}}(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right] + \lambda \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize  $\boldsymbol{\theta}$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

simultaneous update  
for  $j = 0 \dots d$

Use the natural logarithm ( $\ln = \log_e$ ) to cancel with the  $\exp()$  in  $h_{\boldsymbol{\theta}}(\mathbf{x})$

# Gradient Descent for Logistic Regression

$$J_{\text{reg}}(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right] + \lambda \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize  $\boldsymbol{\theta}$
- Repeat until convergence (simultaneous update for  $j = 0 \dots d$ )

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \left[ \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} - \frac{\lambda}{n} \theta_j \right]$$

# Gradient Descent for Logistic Regression

- Initialize  $\theta$
- Repeat until convergence (simultaneous update for  $j = 0 \dots d$ )

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^n \left( h_{\theta} \left( \mathbf{x}^{(i)} \right) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \left[ \sum_{i=1}^n \left( h_{\theta} \left( \mathbf{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)} - \frac{\lambda}{n} \theta_j \right]$$

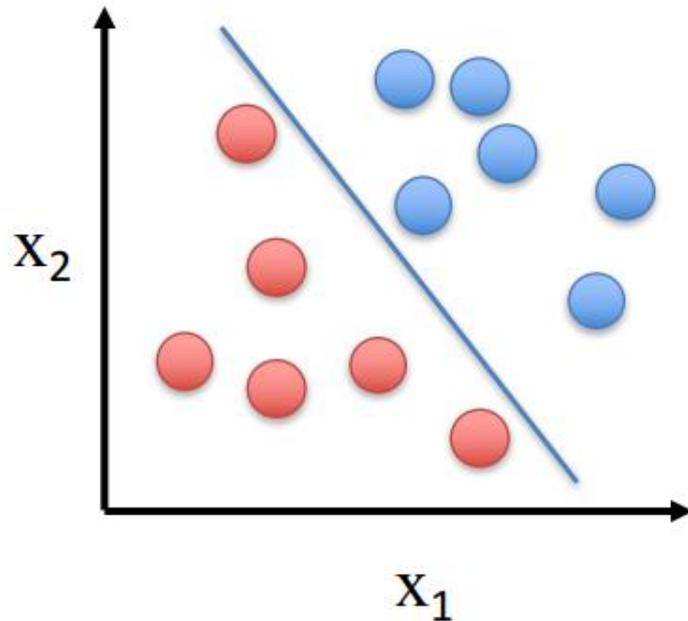
This looks IDENTICAL to linear regression!!!

- Ignoring the  $1/n$  constant
- However, the form of the model is very different:

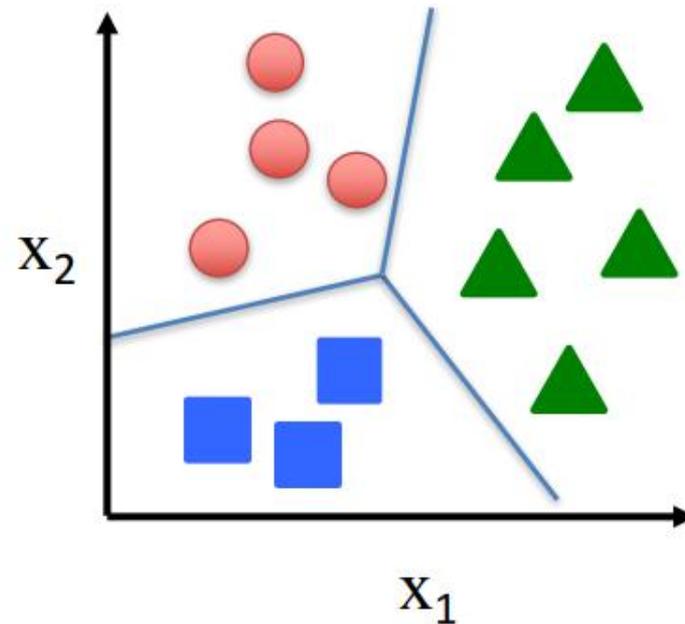
$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$

# Multi-Class Classification

Binary classification:



Multi-class classification:



Disease diagnosis: healthy / cold / flu / pneumonia

Object classification: desk / chair / monitor / bookcase

# Multi-Class Logistic Regression

- For 2 classes:

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})} = \frac{\exp(\theta^T \mathbf{x})}{\boxed{1} + \boxed{\exp(\theta^T \mathbf{x})}}$$

weight assigned to  $y = 0$                       weight assigned to  $y = 1$

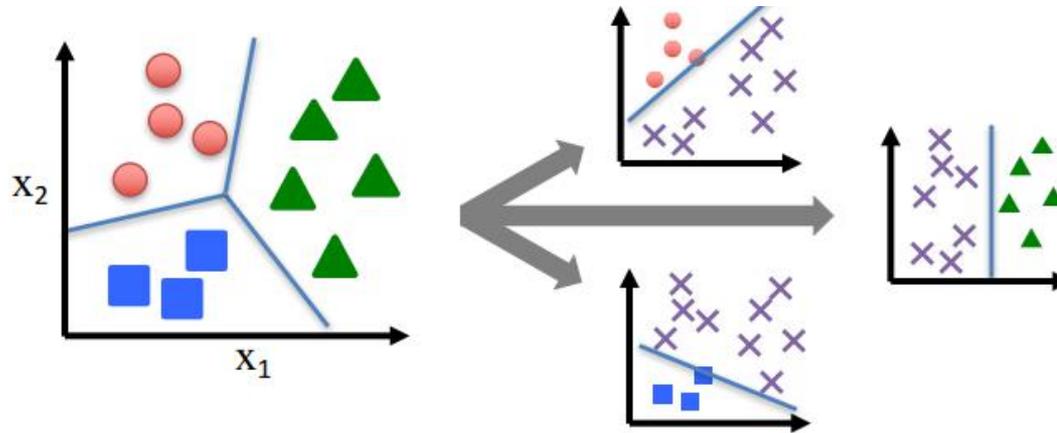
- For  $C$  classes  $\{1, \dots, C\}$ :

$$p(y = c \mid \mathbf{x}; \theta_1, \dots, \theta_C) = \frac{\exp(\theta_c^T \mathbf{x})}{\sum_{c=1}^C \exp(\theta_c^T \mathbf{x})}$$

– Called the **softmax** function

# Multi-Class Logistic Regression

- Split into One vs Rest:



- Train a logistic regression classifier for each class  $i$  to predict the probability that  $y = i$  with

$$h_c(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_c^T \mathbf{x})}{\sum_{c=1}^C \exp(\boldsymbol{\theta}_c^T \mathbf{x})}$$

# Implementing Multi-Class Logistic Regression

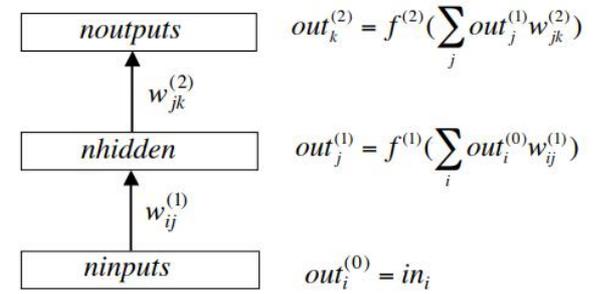
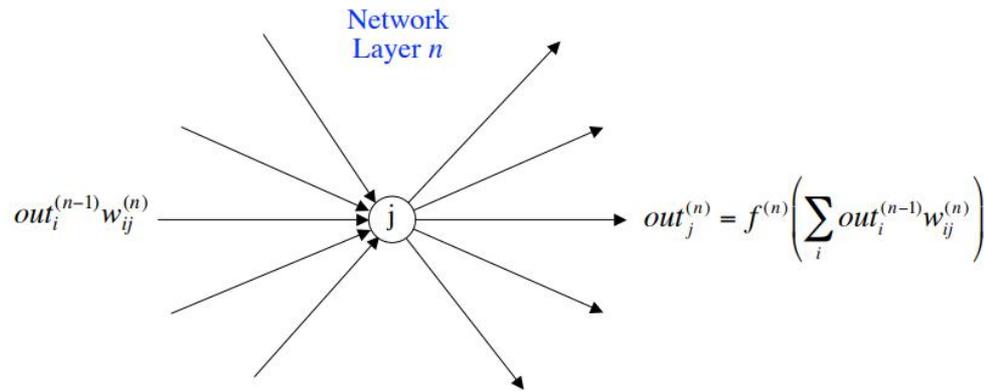
- Use  $h_c(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_c^\top \mathbf{x})}{\sum_{c=1}^C \exp(\boldsymbol{\theta}_c^\top \mathbf{x})}$  as the model for class  $c$
- Gradient descent simultaneously updates all parameters for all models
  - Same derivative as before, just with the above  $h_c(\mathbf{x})$
- Predict class label as the most probable label

$$\max_c h_c(\mathbf{x})$$

# Outline

- Regression Overview
- Linear Regression
- Support Vector Regression
- Logistic Regression
- **Deep Neural Network for Regression**

# DNN Regression



- For a two-layer MLP:

$$out_k^{(2)} = f^{(2)}\left(\sum_j out_j^{(1)} \cdot w_{jk}^{(2)}\right) = f^{(2)}\left(\sum_j f^{(1)}\left(\sum_i in_i w_{ij}^{(1)}\right) \cdot w_{jk}^{(2)}\right)$$

- The network weights are adjusted to minimize an output cost function

$$E_{SSE} = \frac{1}{2} \sum_p \sum_j (targ_j^p - out_j^{(N)p})^2$$

# Computing the Partial Derivatives for Regression

- We use SSE and for a two layer network the linear final outputs can be written:

$$out_k^{(2)} = \sum_j out_j^{(1)} \cdot w_{jk}^{(2)} = \sum_j f\left(\sum_i in_i w_{ij}^{(1)}\right) \cdot w_{jk}^{(2)}$$

- We can then use the chain rules for derivatives, as for the Single Layer Perceptron, to give the derivatives with respect to the two sets of weights

$$\frac{\partial E_{SSE}(\{w_{ij}^{(n)}\})}{\partial w_{hl}^{(m)}} = - \sum_p \sum_k (targ_k - out_k^{(2)}) \cdot \frac{\partial out_k^{(2)}}{\partial w_{hl}^{(m)}}$$

$$\frac{\partial out_k^{(2)}}{\partial w_{hl}^{(2)}} = \sum_j out_j^{(1)} \frac{\partial w_{jk}^{(2)}}{\partial w_{hl}^{(2)}} = \sum_j out_j^{(1)} \cdot \delta_{jh} \cdot \delta_{kl} = out_h^{(1)} \cdot \delta_{kl}$$

$$\frac{\partial out_k^{(2)}}{\partial w_{hl}^{(1)}} = \sum_j f'\left(\sum_i in_i w_{ij}^{(1)}\right) \cdot \left(\sum_m in_m \frac{\partial w_{mj}^{(1)}}{\partial w_{hl}^{(1)}}\right) \cdot w_{jk}^{(2)} = f'\left(\sum_i in_i w_{il}^{(1)}\right) in_h \cdot w_{lk}^{(2)}$$

# Deriving the Back Propagation Algorithm for Regression

- All we now have to do is substitute our derivatives into the weight update equations

$$\Delta w_{hl}^{(2)} = \eta \sum_p \sum_k (targ_k - out_k^{(2)}) \cdot out_h^{(1)} \cdot \delta_{kl} = \eta \sum_p (targ_l - out_l^{(2)}) out_h^{(1)}$$

$$\Delta w_{hl}^{(1)} = \eta \sum_p \sum_k (targ_k - out_k^{(2)}) \cdot f' \left( \sum_i in_i w_{il}^{(1)} \right) \cdot w_{lk}^{(2)} \cdot in_h$$

- Then if the transfer function  $f(x)$  is a Sigmoid we can use  $f'(x) = f(x) \cdot (1 - f(x))$  to give

$$\Delta w_{hl}^{(2)} = \eta \sum_p (targ_l - out_l^{(2)}) out_h^{(1)}$$

$$\Delta w_{hl}^{(1)} = \eta \sum_p \sum_k (targ_k - out_k^{(2)}) w_{lk}^{(2)} \cdot out_l^{(1)} \cdot (1 - out_l^{(1)}) \cdot in_h$$

- These equations constitute the Back-Propagation Learning Algorithm for Regression.

# Classification + Localization: Task

**Classification:** C classes

**Input:** Image

**Output:** Class label

**Evaluation metric:** Accuracy



CAT

**Localization:**

**Input:** Image

**Output:** Box in the image (x, y, w, h)

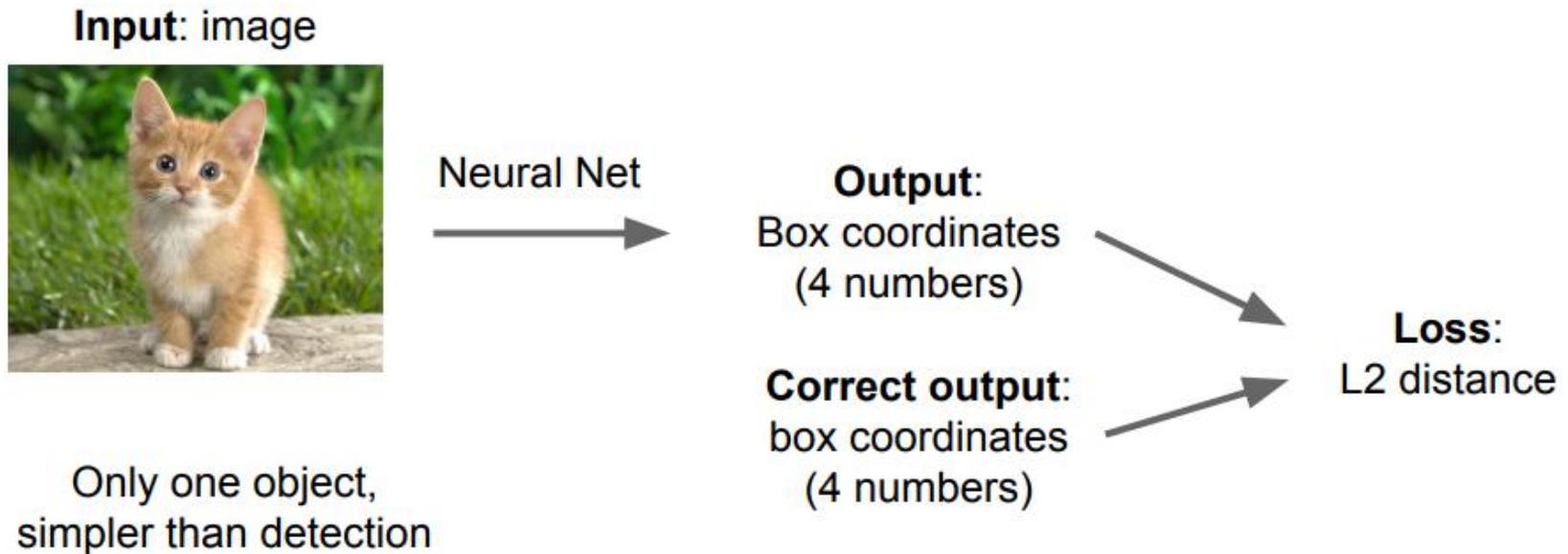
**Evaluation metric:** Intersection over Union



(x, y, w, h)

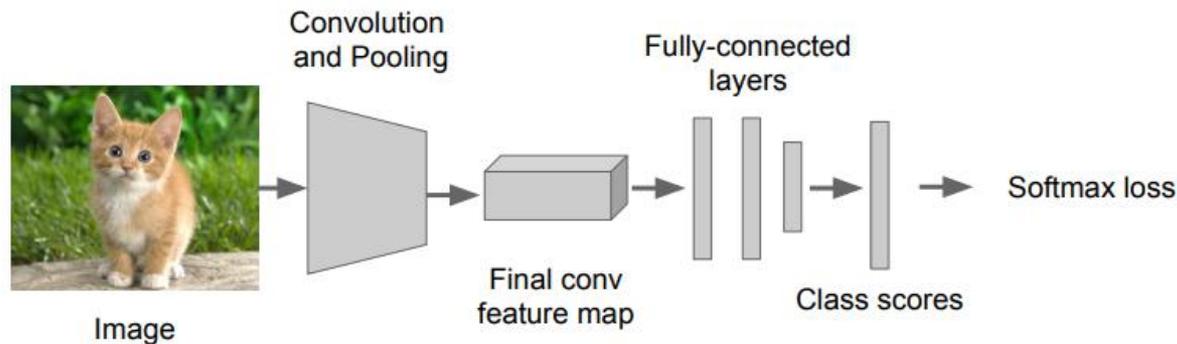
**Classification + Localization:** Do both

# Idea #1: Localization as Regression



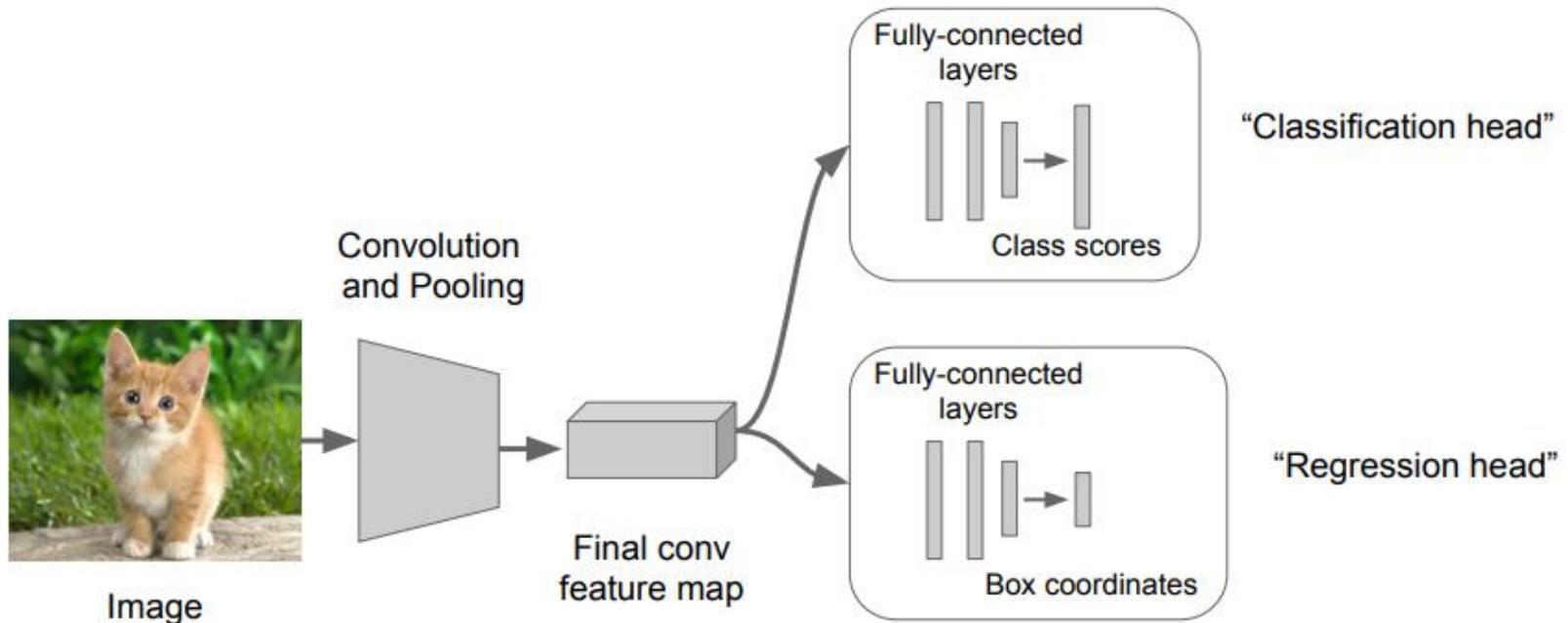
# Simple Recipe for Classification + Localization

- Step 1: Train (or download) a classification model (AlexNet, VGG, GoogLeNet)



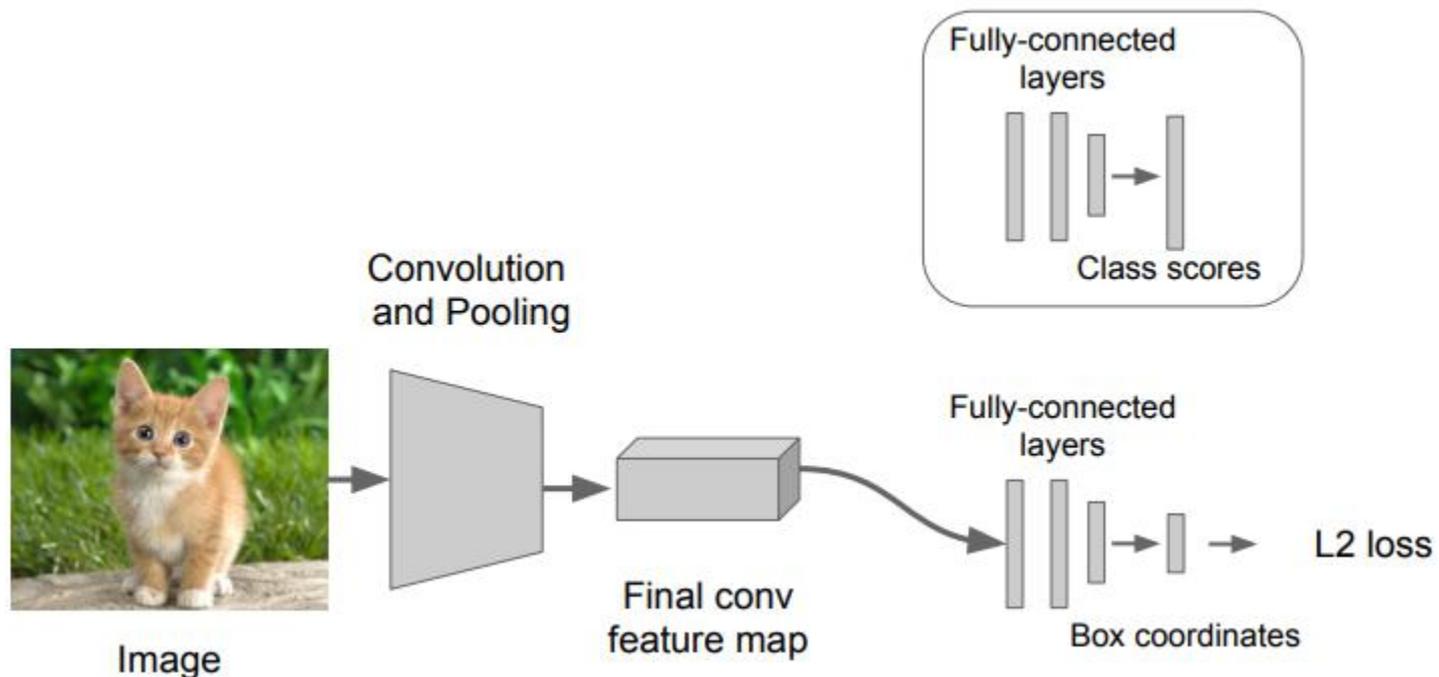
# Simple Recipe for Classification + Localization

- Step 2: Attach new fully-connected “regression head” to the network



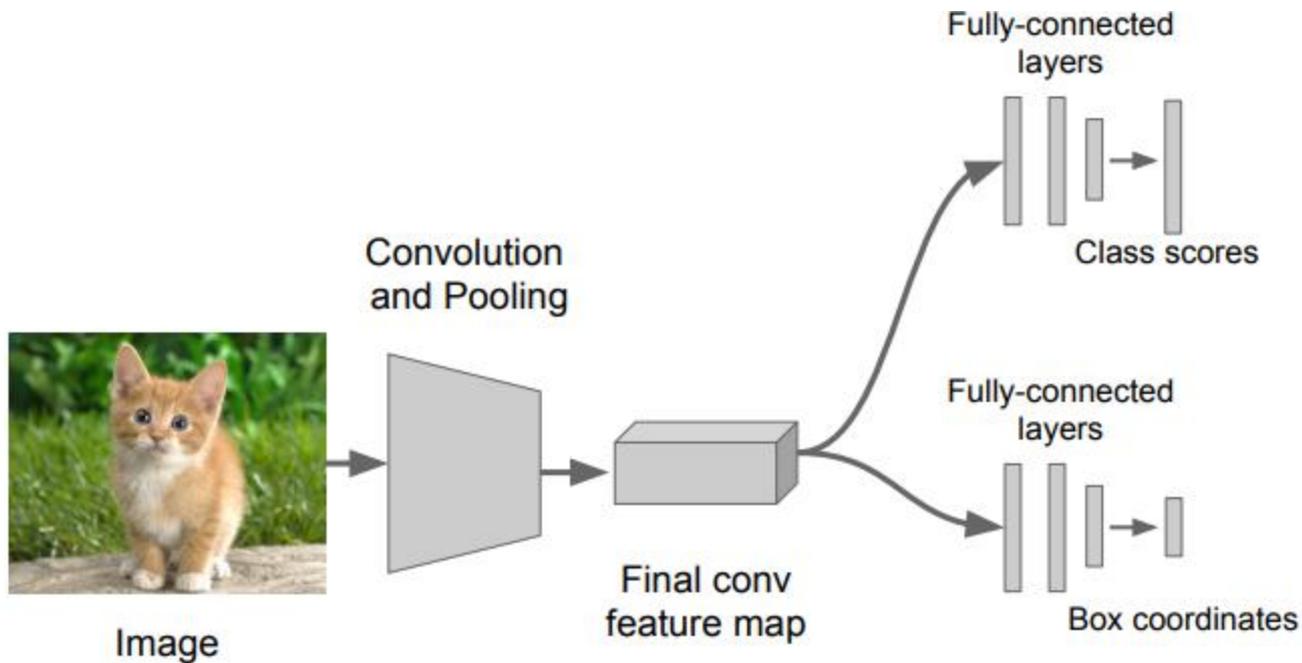
# Simple Recipe for Classification + Localization

- Step 3: Train the regression head only with SGD and L2 loss



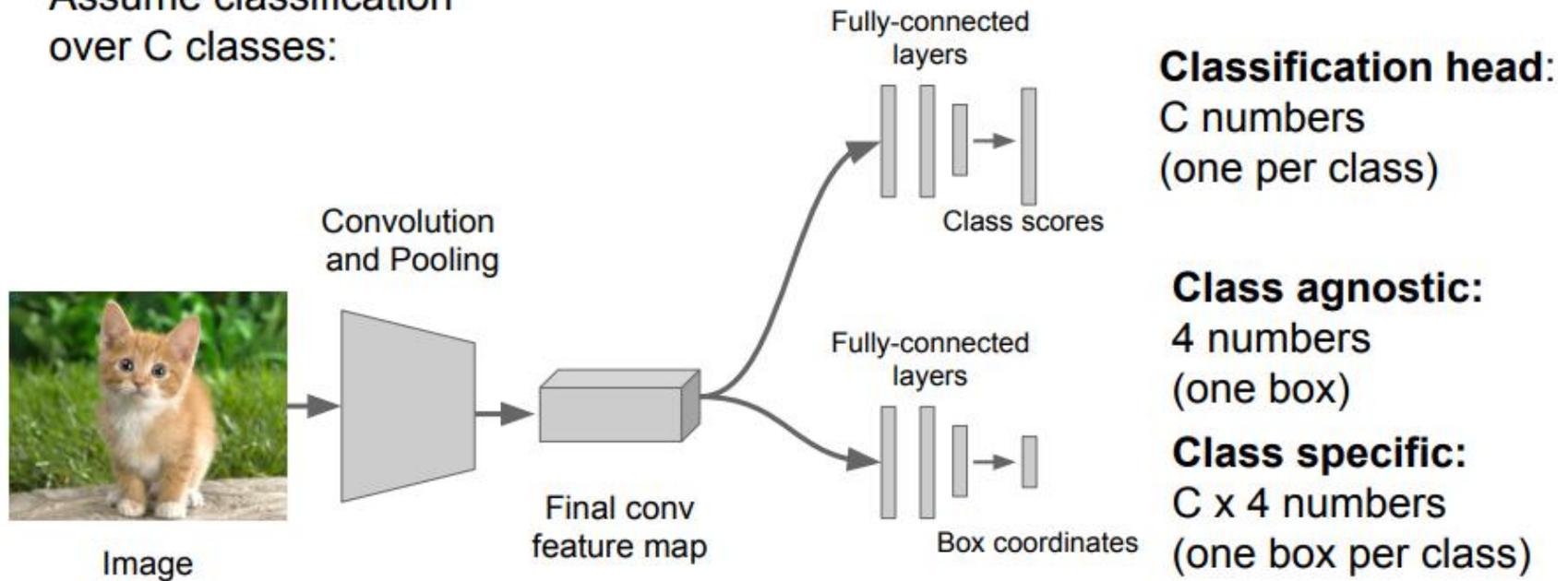
# Simple Recipe for Classification + Localization

- Step 4: At test time use both heads

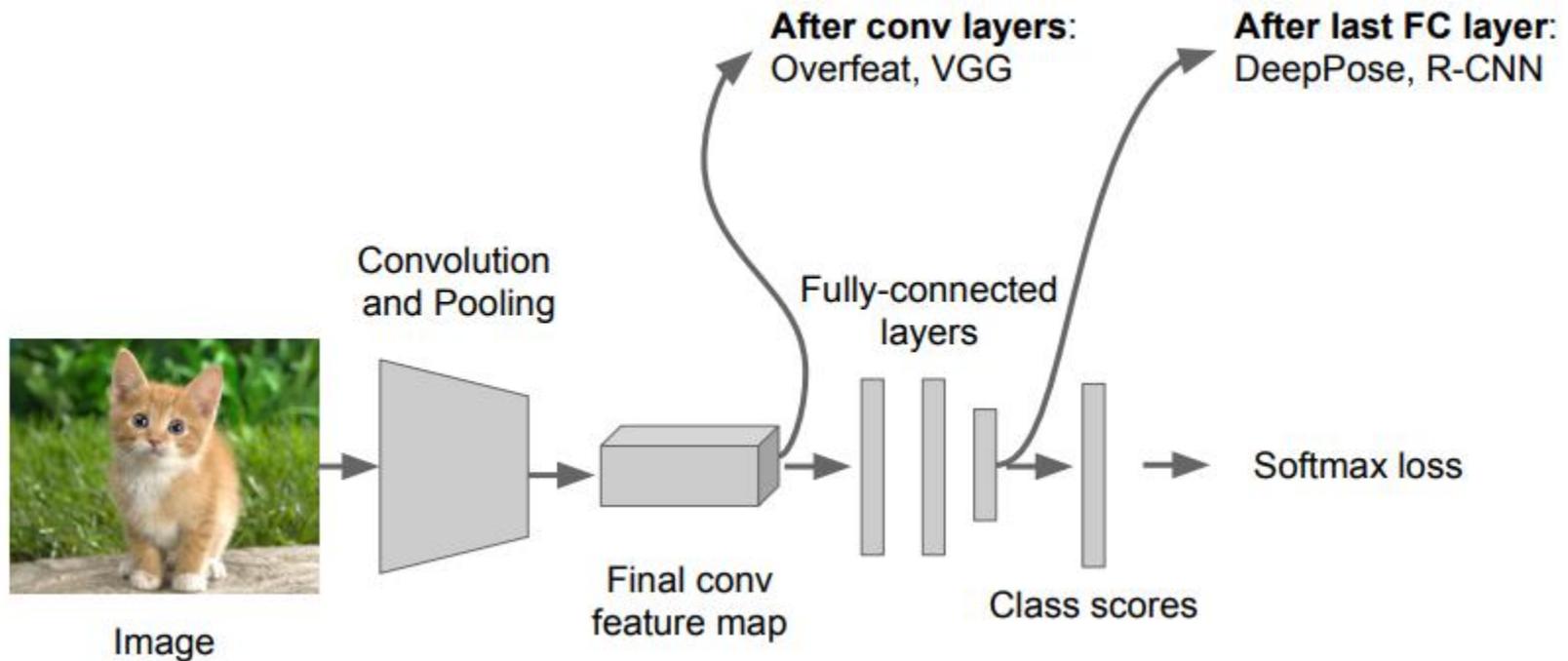


# Per-class vs class agnostic regression

Assume classification over  $C$  classes:



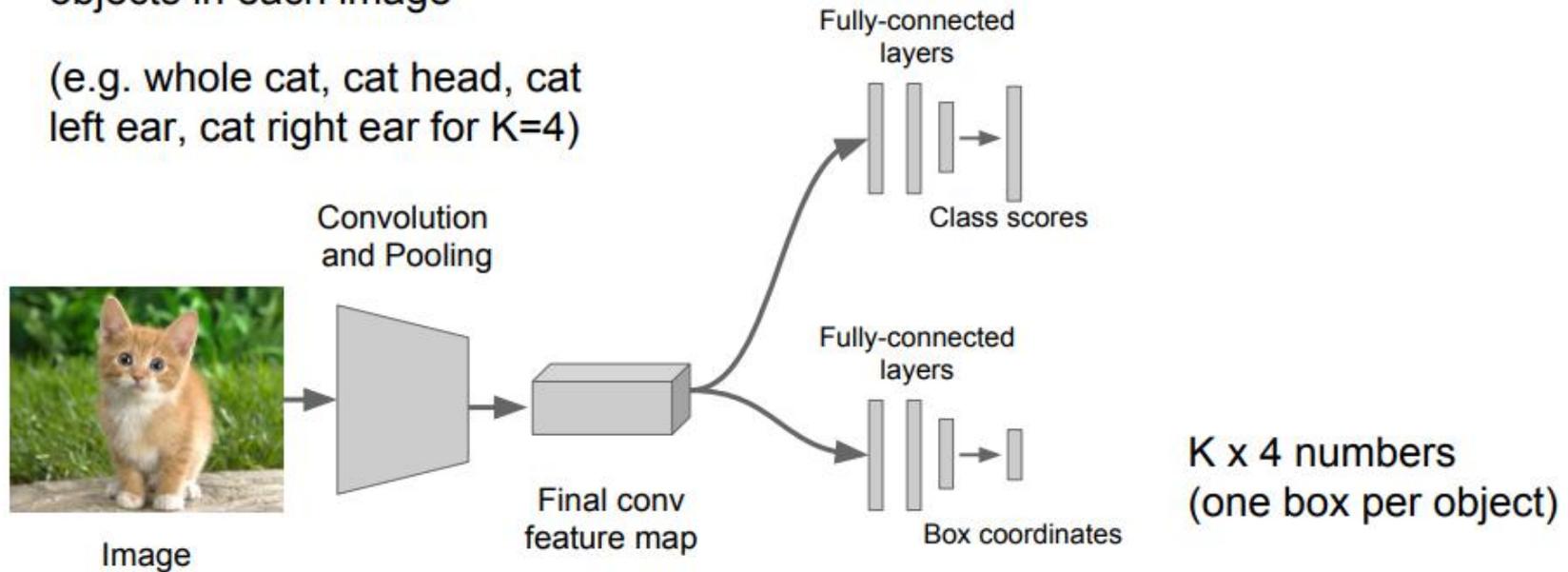
# Where to attach the regression head?



# Aside: Localizing multiple objects

Want to localize **exactly**  $K$  objects in each image

(e.g. whole cat, cat head, cat left ear, cat right ear for  $K=4$ )



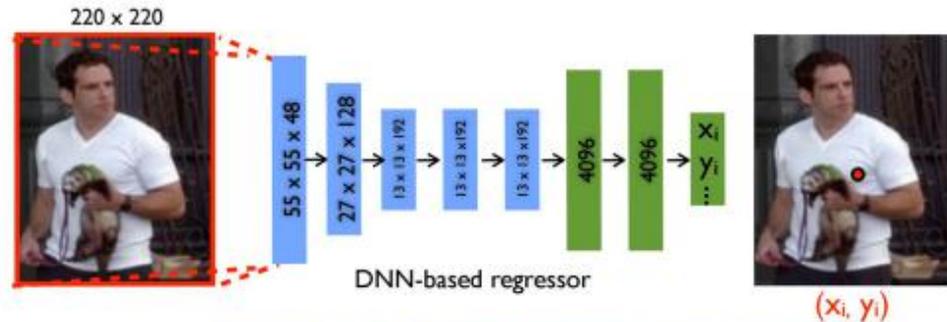
# Aside: Human Pose Estimation

Represent a person by  $K$  joints

Regress  $(x, y)$  for each joint from last fully-connected layer of AlexNet

(Details: Normalized coordinates, iterative refinement)

Toshev and Szegedy, "DeepPose: Human Pose Estimation via Deep Neural Networks", CVPR 2014



# DNN Regression Applications

- Great results in:
  - Computer Vision
    - Object Localization / Detection as DNN Regression
    - Self-driving Steering Command Prediction
    - Human Pose Regression
  - Finance
    - Currency Exchange Rate
    - Stock Price Prediction
    - Forecasting Financial Time Series
    - Crude Oil Price Prediction

# DNN Regression Applications

- **Great results in:**
  - **Atmospheric Sciences**
    - Air Quality Prediction
    - Carbon Dioxide Pollution Prediction
    - Ozone Concentration Modeling
    - Sulphur Dioxide Concentration Prediction
  - **Infrastructure**
    - Road Tunnel Cost Estimation
    - Highway Engineering Cost Estimation
  - **Geology / Physics**
    - Meteorology and Oceanography Application
    - Pacific Sea Surface Temperature Prediction
    - Hydrological Modeling

# *Q & A*