

Lecture 8: 3 February, 2026

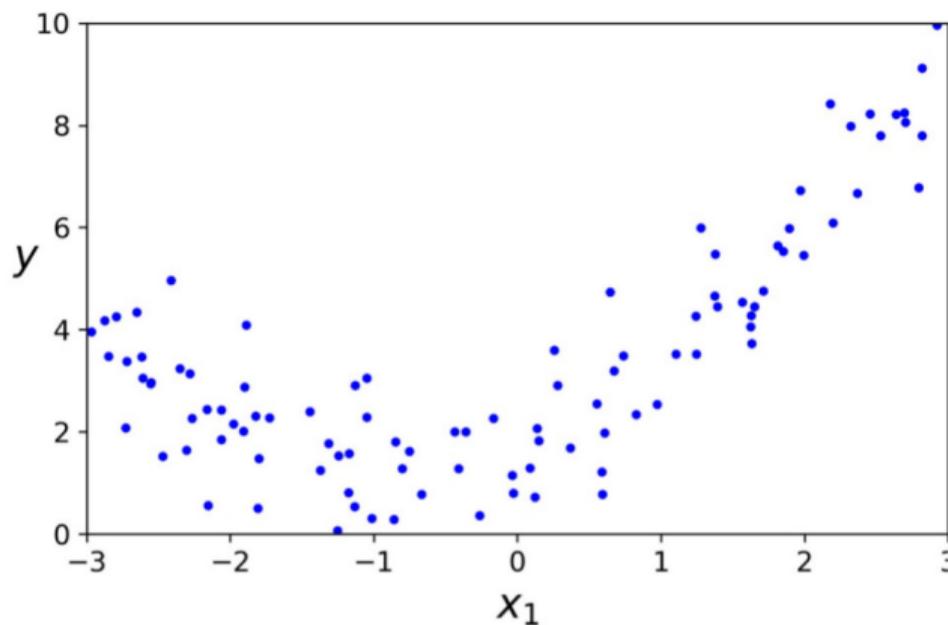
Madhavan Mukund

<https://www.cmi.ac.in/~madhavan>

Data Mining and Machine Learning
January–April 2026

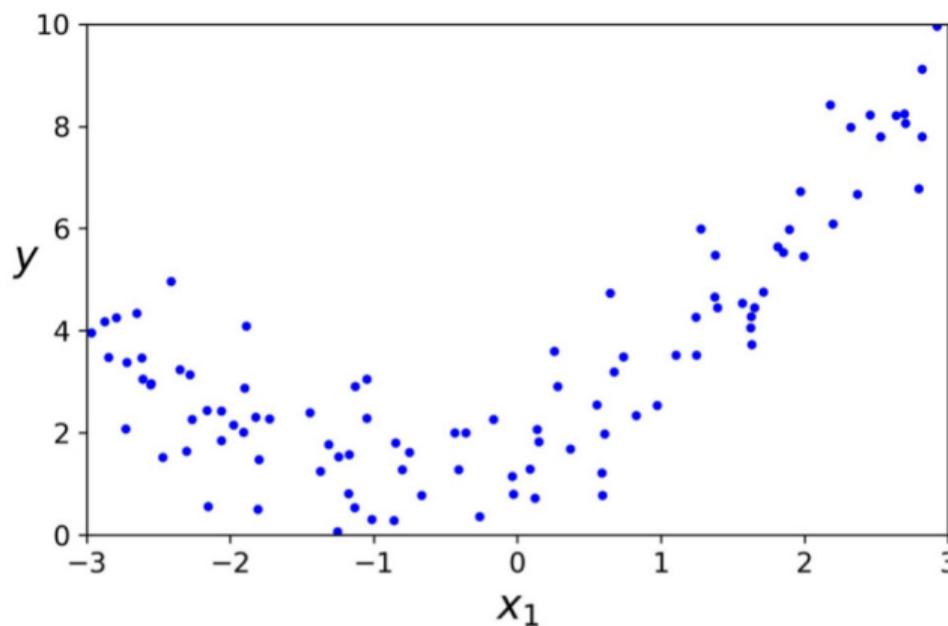
Decision trees for regression

- Can we use decision trees for regression?



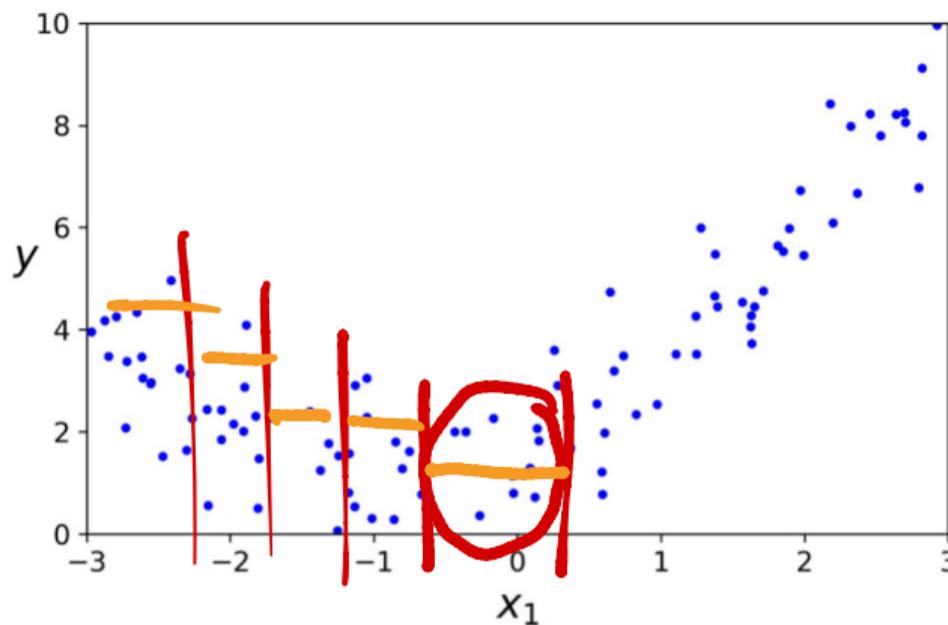
Decision trees for regression

- Can we use decision trees for regression?
- Partition the input into intervals



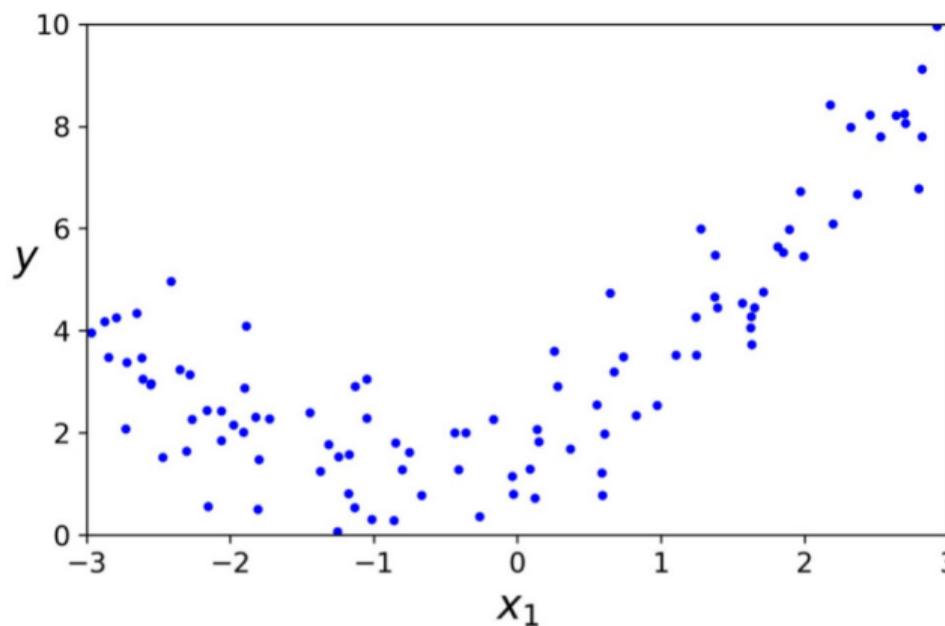
Decision trees for regression

- Can we use decision trees for regression?
- Partition the input into intervals
- For each interval, predict mean value of output, instead of majority class



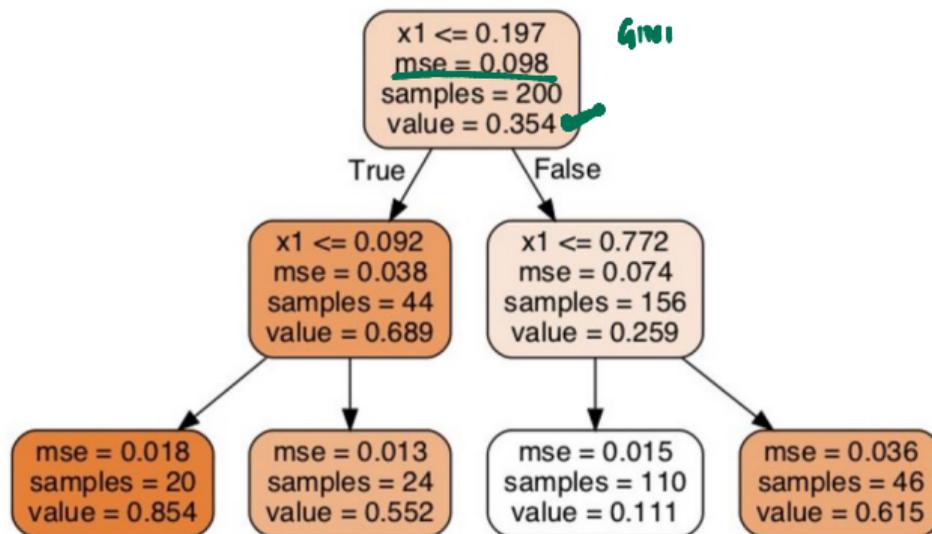
Decision trees for regression

- Can we use decision trees for regression?
- Partition the input into intervals
- For each interval, predict mean value of output, instead of majority class
- Regression tree



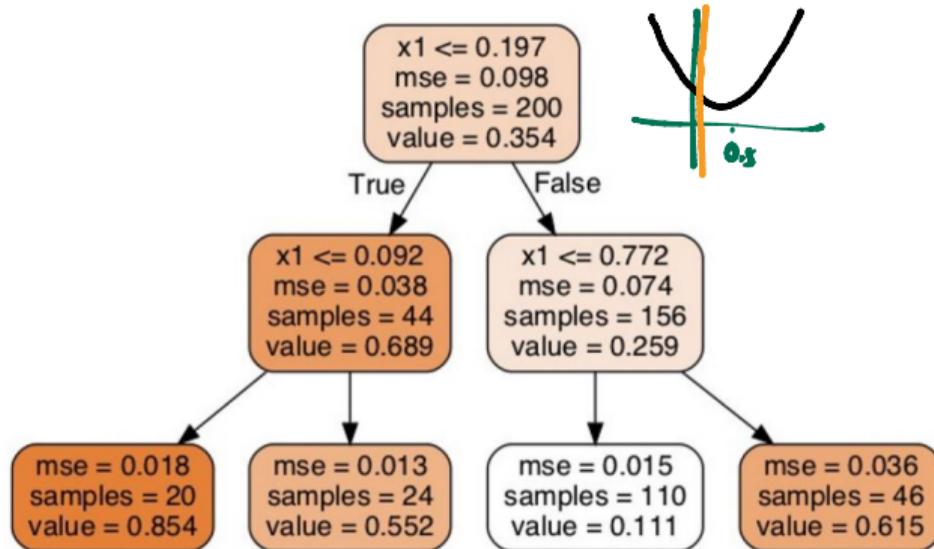
Decision trees for regression

- Regression tree for noisy quadratic centered around $x_1 = 0.5$



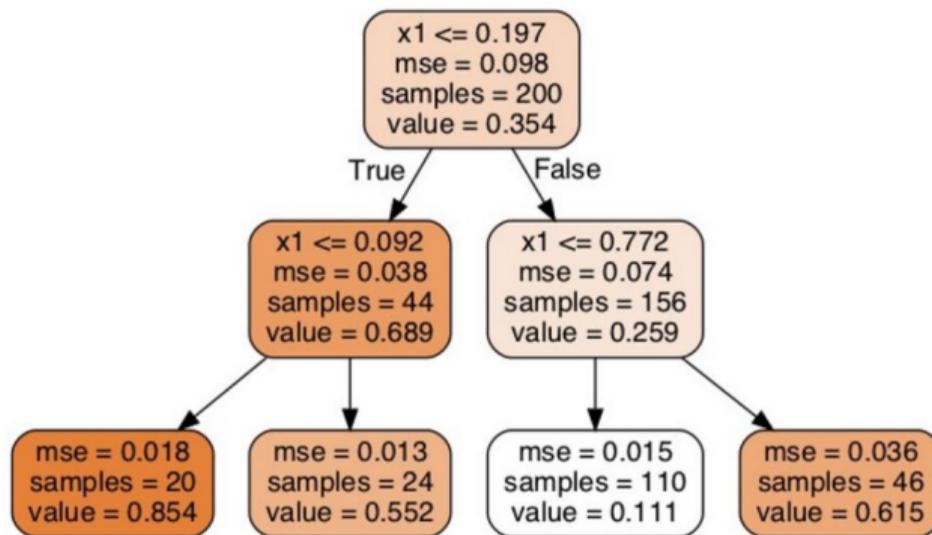
Decision trees for regression

- Regression tree for noisy quadratic centered around $x_1 = 0.5$
- For each node, the output is the mean y value for the current set of points



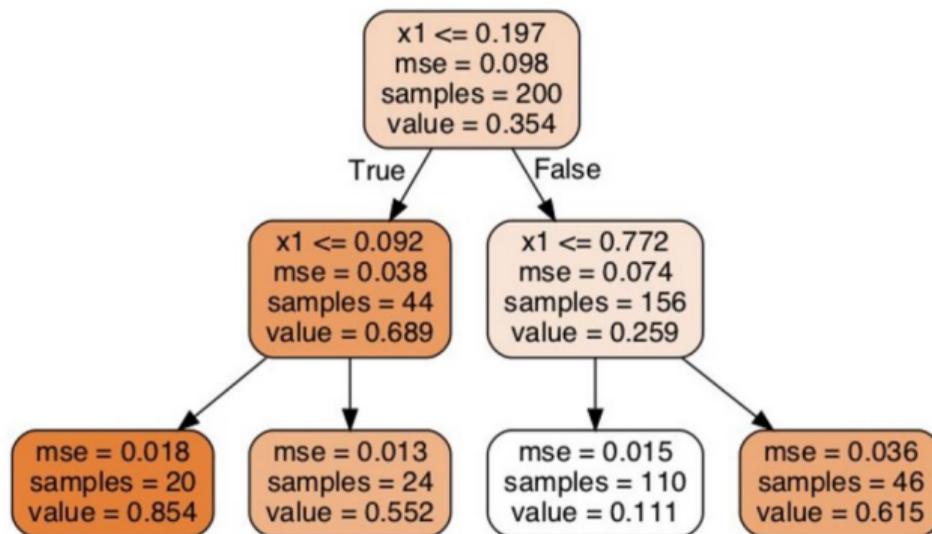
Decision trees for regression

- Regression tree for noisy quadratic centered around $x_1 = 0.5$
- For each node, the output is the mean y value for the current set of points
- Instead of impurity, use mean squared error (MSE) as cost function



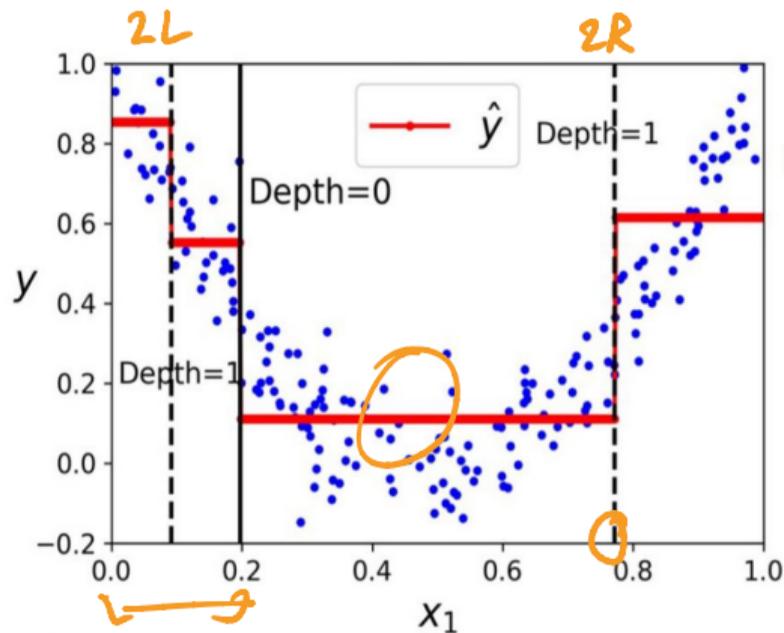
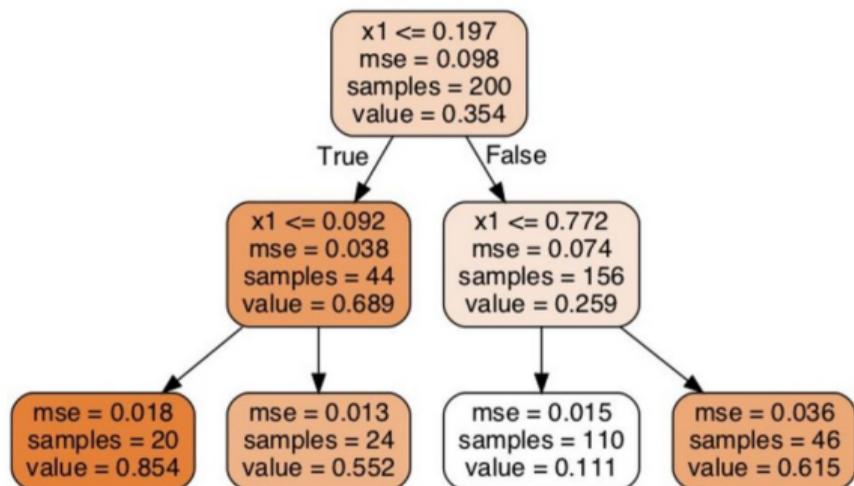
Decision trees for regression

- Regression tree for noisy quadratic centered around $x_1 = 0.5$
- For each node, the output is the mean y value for the current set of points
- Instead of impurity, use mean squared error (MSE) as cost function
- Choose a split that minimizes MSE



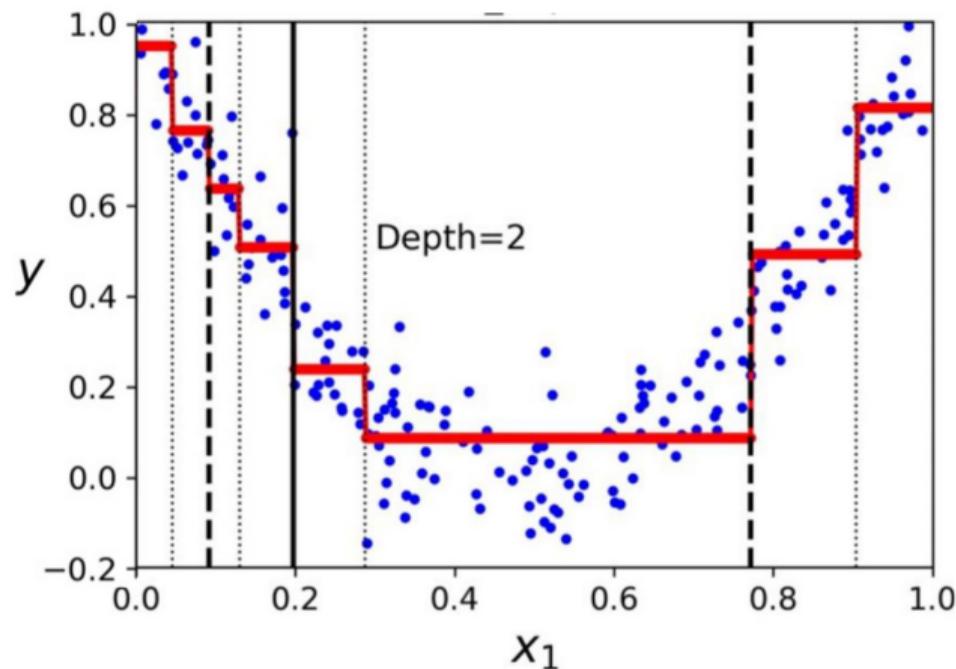
Regression trees

■ Approximation using regression tree



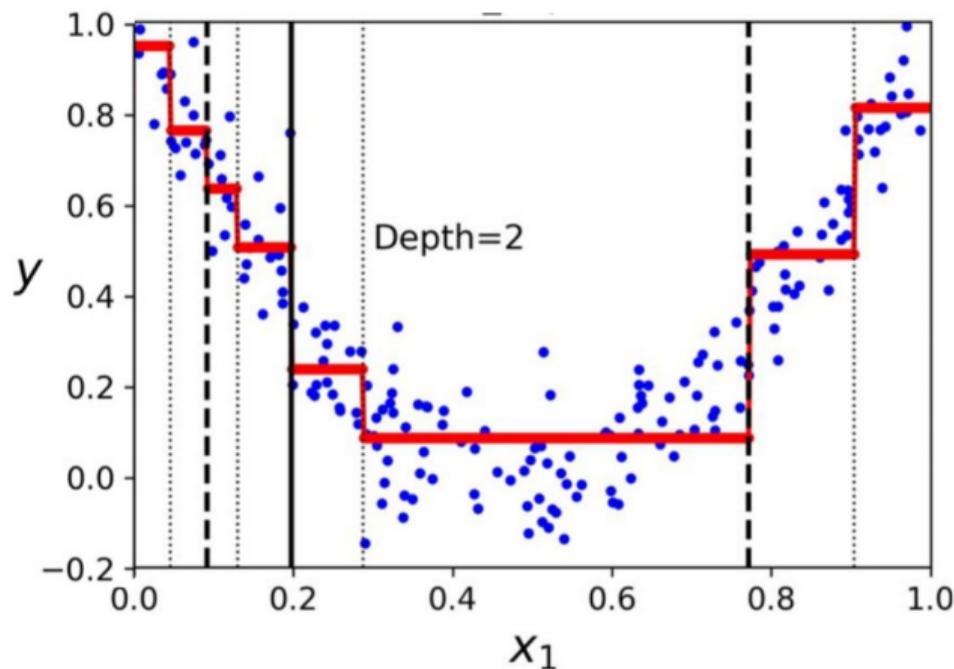
Regression trees

- Extend the regression tree one more level to get a finer approximation



Regression trees

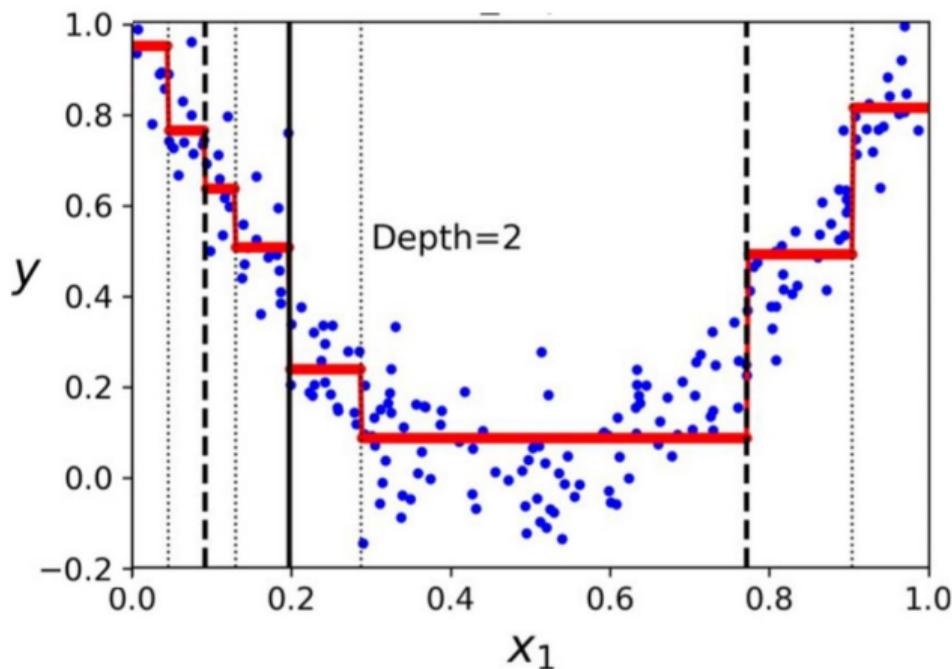
- Extend the regression tree one more level to get a finer approximation
- Set a threshold on MSE to decide when to stop



Regression trees

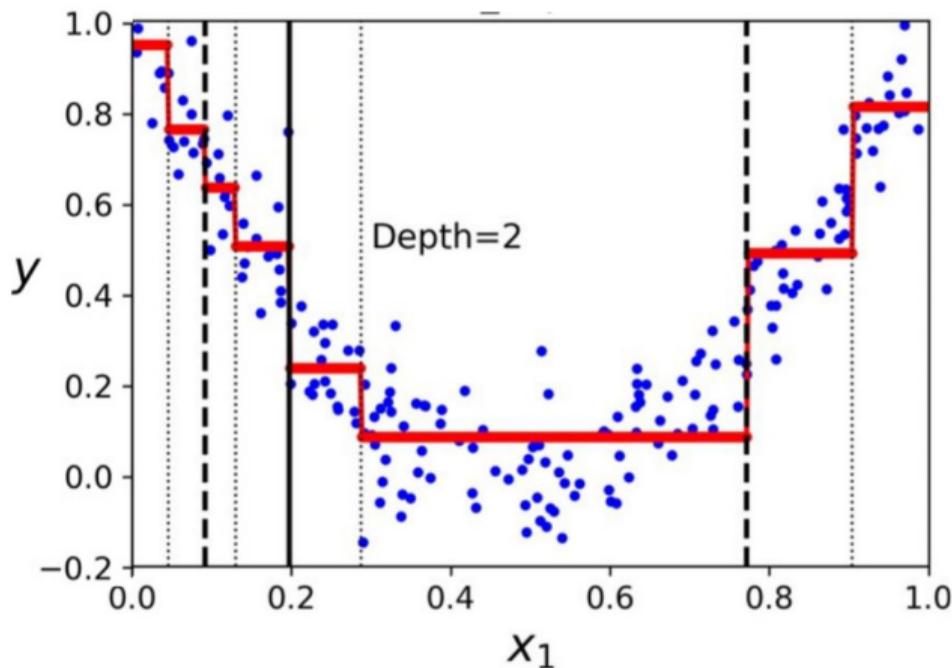
- Extend the regression tree one more level to get a finer approximation
- Set a threshold on MSE to decide when to stop
- Classification and Regression Trees (CART)

Leo Breiman



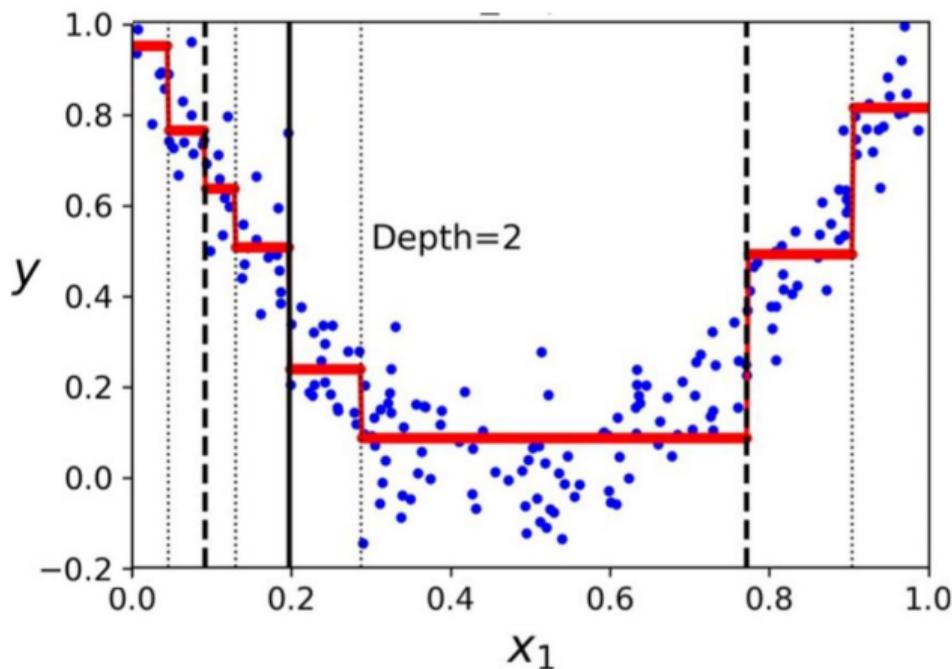
Regression trees

- Extend the regression tree one more level to get a finer approximation
- Set a threshold on MSE to decide when to stop
- Classification and Regression Trees (CART)
 - Combined algorithm for both use cases



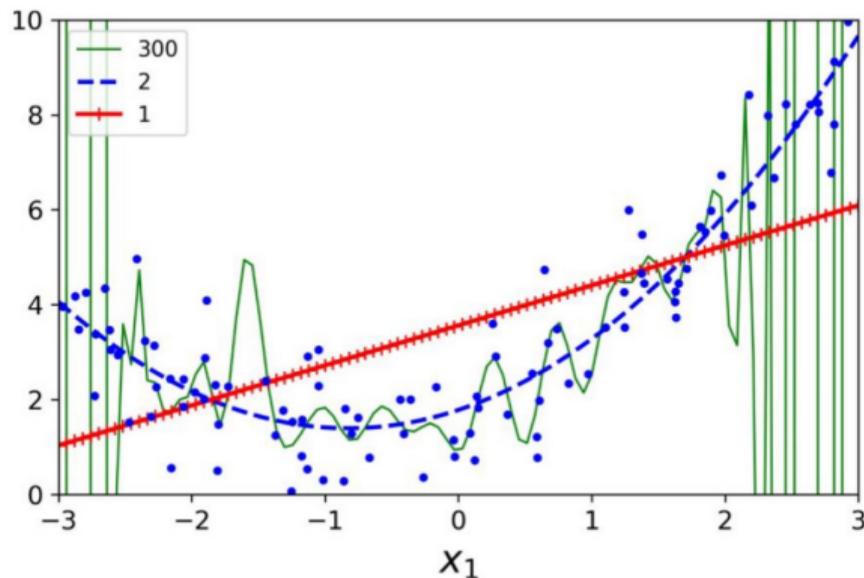
Regression trees

- Extend the regression tree one more level to get a finer approximation
- Set a threshold on MSE to decide when to stop
- **Classification and Regression Trees (CART)**
 - Combined algorithm for both use cases
- Programming libraries typically provide CART implementation



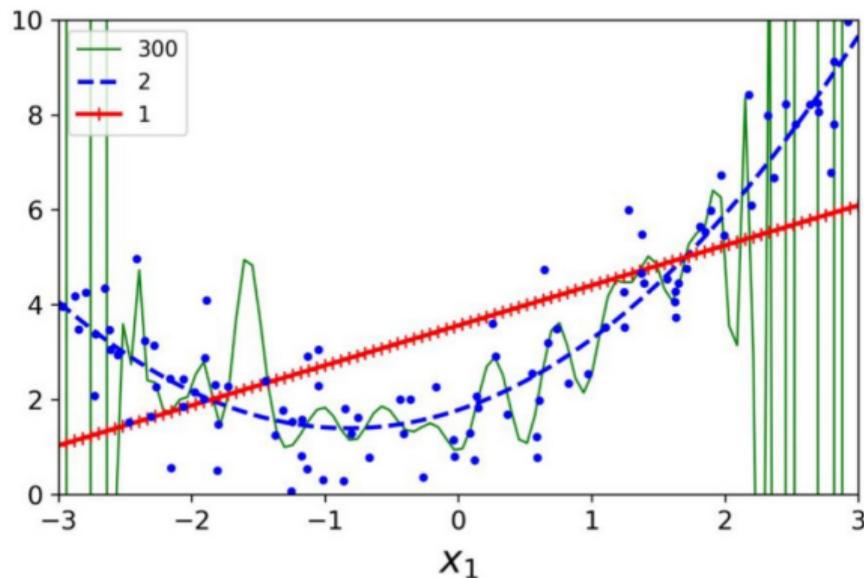
Overfitting

- Overfitting: model too specific to training data, does not generalize well



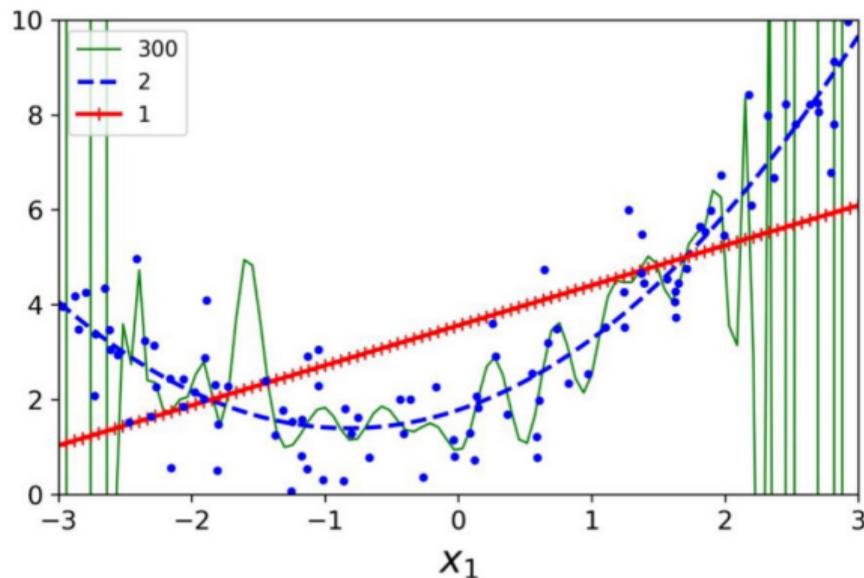
Overfitting

- Overfitting: model too specific to training data, does not generalize well
- Regression — use regularization to penalize model complexity



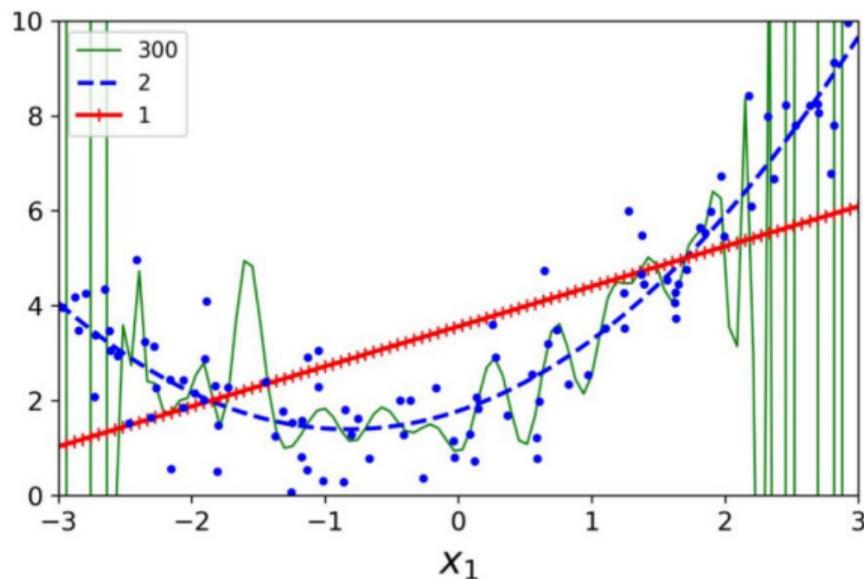
Overfitting

- Overfitting: model too specific to training data, does not generalize well
- Regression — use regularization to penalize model complexity
- What about decision trees?



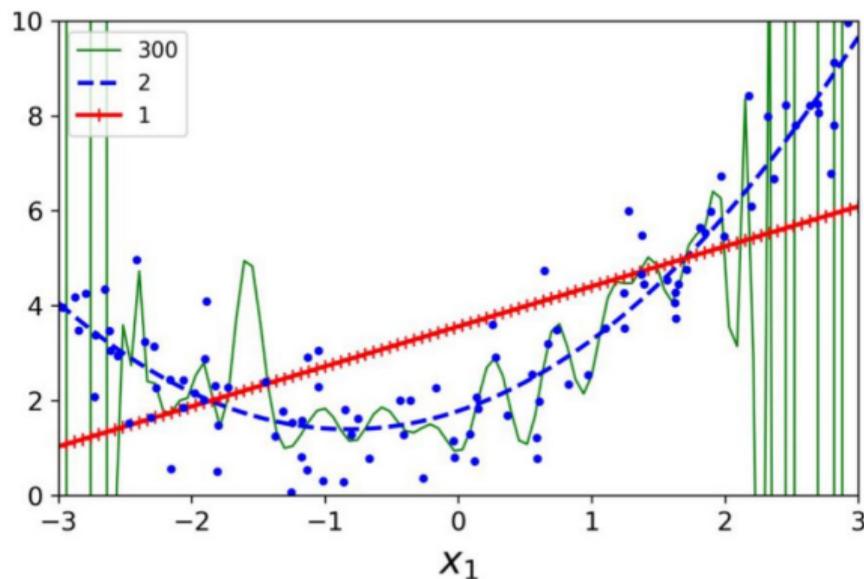
Overfitting

- Overfitting: model too specific to training data, does not generalize well
- Regression — use regularization to penalize model complexity
- What about decision trees?
- Deep, complex trees ask too many questions



Overfitting

- Overfitting: model too specific to training data, does not generalize well
- Regression — use regularization to penalize model complexity
- What about decision trees?
- Deep, complex trees ask too many questions
- Prefer shallow, simple trees



Tree pruning

- Remove leaves to improve generalization

Tree pruning

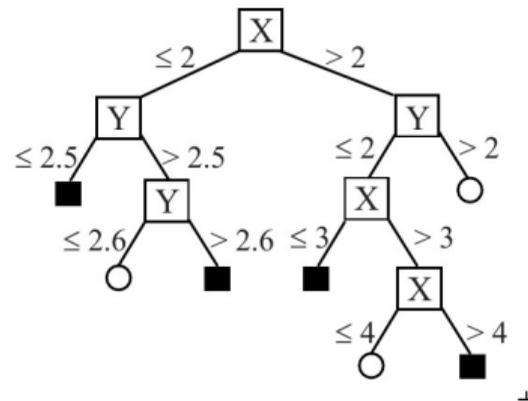
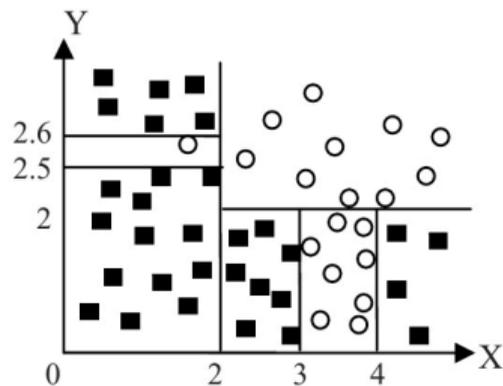
- Remove leaves to improve generalization
- Top-down pruning
 - Fix a maximum depth when building the tree
 - How to decide the depth in advance?

Tree pruning

- Remove leaves to improve generalization
- Top-down pruning
 - Fix a maximum depth when building the tree
 - How to decide the depth in advance?
- Bottom-up pruning
 - Build the full tree
 - Remove a leaf if the reduced tree generalizes better
 - How do we measure this?

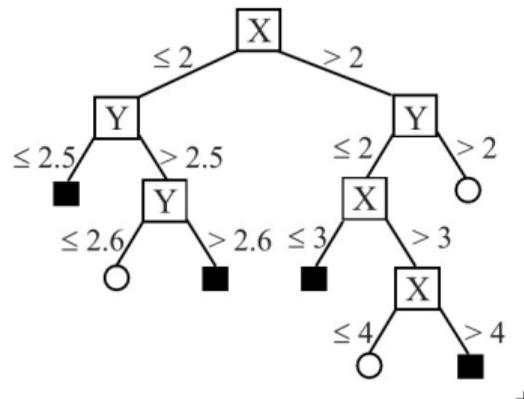
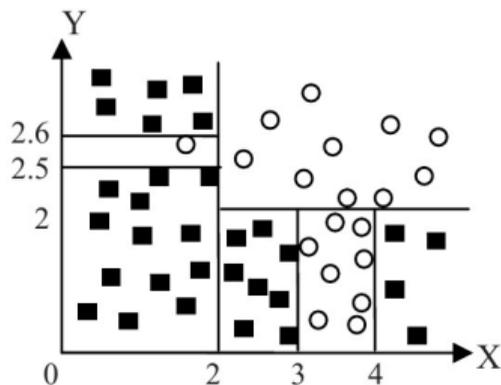
Tree pruning

Overfitted tree

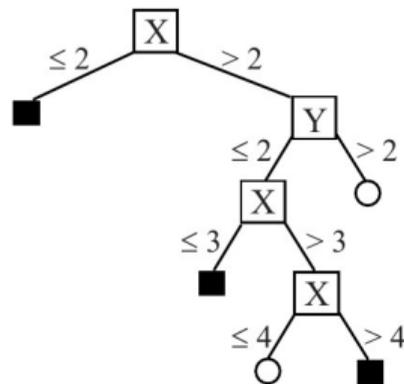
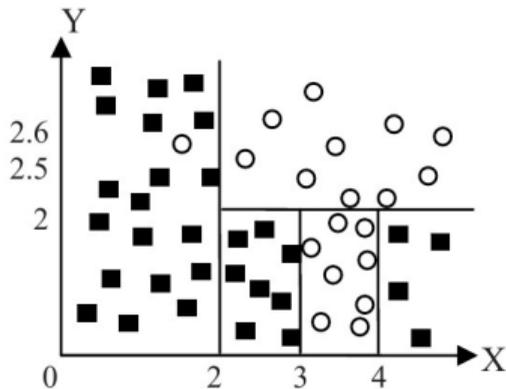


Tree pruning

Overfitted tree



Pruned tree



Bottom up tree pruning

- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?

Bottom up tree pruning

- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?
- Check performance on a test set

Bottom up tree pruning

- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?
- Check performance on a test set
- Use **sampling theory** [Quinlan]

Entropy

Bottom up tree pruning

- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?
- Check performance on a test set
- Use **sampling theory** [Quinlan]
- Given n coin tosses with h heads, estimate probability of heads as h/n

Bottom up tree pruning

- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?
- Check performance on a test set
- Use **sampling theory** [Quinlan]
- Given n coin tosses with h heads, estimate probability of heads as h/n
 - Estimate comes with a confidence interval: $h/n \pm \delta$

Bottom up tree pruning

- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?
- Check performance on a test set
- Use **sampling theory** [Quinlan]
- Given n coin tosses with h heads, estimate probability of heads as h/n
 - Estimate comes with a confidence interval: $h/n \pm \delta$
 - As n increases, δ reduces: **7 heads out of 10** vs **70 out of 100** vs **700 out of 1000**

Bottom up tree pruning

- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?
- Check performance on a test set
- Use **sampling theory** [Quinlan]
- Given n coin tosses with h heads, estimate probability of heads as h/n
 - Estimate comes with a confidence interval: $h/n \pm \delta$
 - As n increases, δ reduces: **7 heads out of 10** vs **70 out of 100** vs **700 out of 1000**
- Impure node, majority prediction, compute confidence interval

Bottom up tree pruning

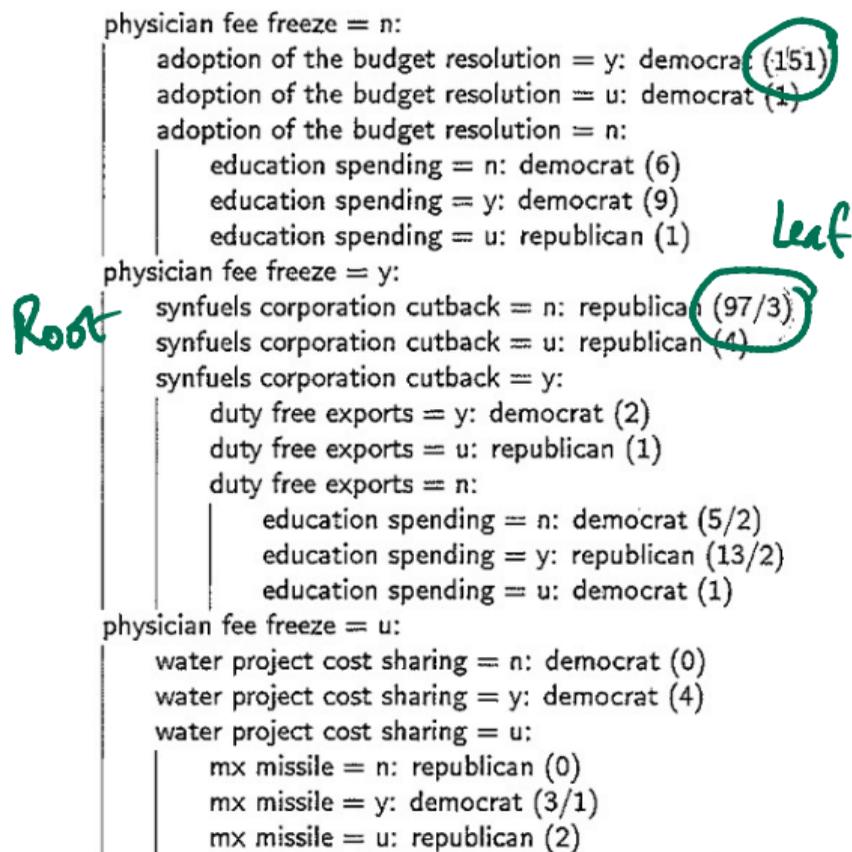
- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?
- Check performance on a test set
- Use **sampling theory** [Quinlan]
- Given n coin tosses with h heads, estimate probability of heads as h/n
 - Estimate comes with a confidence interval: $h/n \pm \delta$
 - As n increases, δ reduces: **7 heads out of 10** vs **70 out of 100** vs **700 out of 1000**
- Impure node, majority prediction, compute confidence interval
- Pruning leaves creates a larger impure sample one level above

Bottom up tree pruning

- Build the full tree, remove leaf if the reduced tree generalizes better
- How do we measure this?
- Check performance on a test set
- Use **sampling theory** [Quinlan]
- Given n coin tosses with h heads, estimate probability of heads as h/n
 - Estimate comes with a confidence interval: $h/n \pm \delta$
 - As n increases, δ reduces: **7 heads out of 10** vs **70 out of 100** vs **700 out of 1000**
- Impure node, majority prediction, compute confidence interval
- Pruning leaves creates a larger impure sample one level above
- Does the confidence interval decrease (improve)?

Example: Predict party from voting pattern [Quinlan]

- Predict party affiliation of US legislators based on voting pattern
 - Read the tree from left to right



Example: Predict party from voting pattern [Quinlan]

- Predict party affiliation of US legislators based on voting pattern
 - Read the tree from left to right
- After pruning, drastically simplified tree

Confidence



physician fee freeze = n: democrat (168/2.6)
physician fee freeze = y: republican (123/13.9)
physician fee freeze = u:
| mx missile = n: democrat (3/1.1)
| mx missile = y: democrat (4/2.2)
| mx missile = u: republican (2/1)

Example: Predict party from voting pattern [Quinlan]

- Predict party affiliation of US legislators based on voting pattern
 - Read the tree from left to right
- After pruning, drastically simplified tree
- Quinlan's comment on his use of sampling theory for post-pruning

Now, this description does violence to statistical notions of sampling and confidence limits, so the reasoning should be taken with a large grain of salt. Like many heuristics with questionable underpinnings, however, the estimates it produces seem frequently to yield acceptable results.

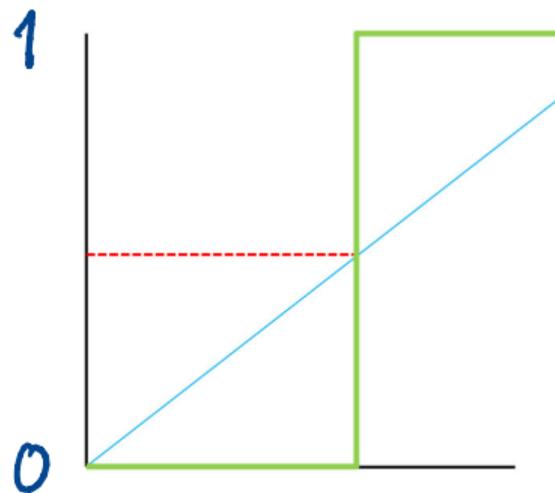
```
physician fee freeze = n: democrat (168/2.6)
physician fee freeze = y: republican (123/13.9)
physician fee freeze = u:
|
| mx missile = n: democrat (3/1.1)
| mx missile = y: democrat (4/2.2)
| mx missile = u: republican (2/1)
```

Regression for classification

- Regression line
- Set a threshold
- Classifier
 - Output below threshold : 0 (No)
 - Output above threshold : 1 (Yes)

Regression for classification

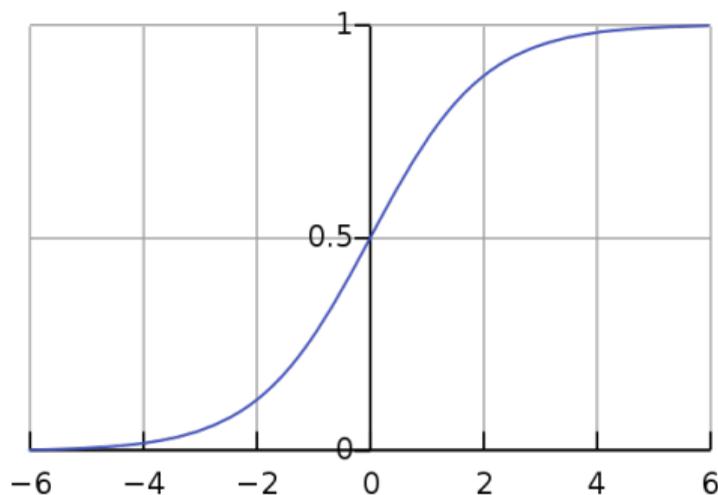
- Regression line
- Set a threshold
- Classifier
 - Output below threshold : 0 (No)
 - Output above threshold : 1 (Yes)
- Classifier output is a step function



Smoothen the step

- Sigmoid function

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



Smoothen the step

- Sigmoid function

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

- Input z is output of our regression

$$\sigma(z) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \dots + \theta_k x_k)}}$$

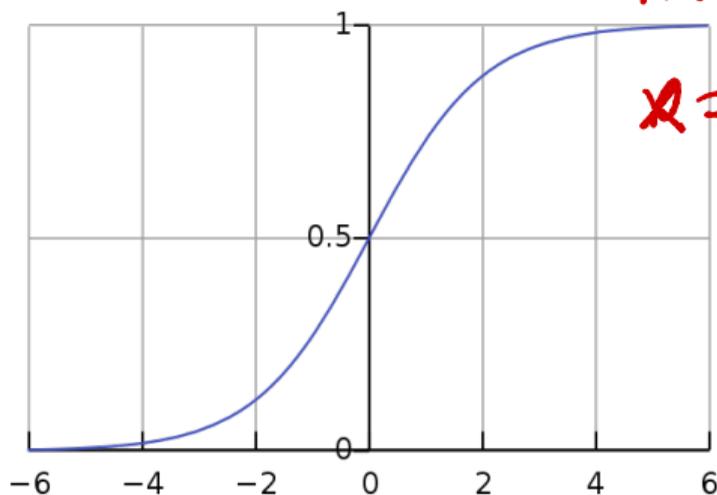
linear output = z

$$z = ax + b$$

At $z=0$

$$ax + b = 0$$

$$x = -\frac{b}{a}$$



Smoothen the step

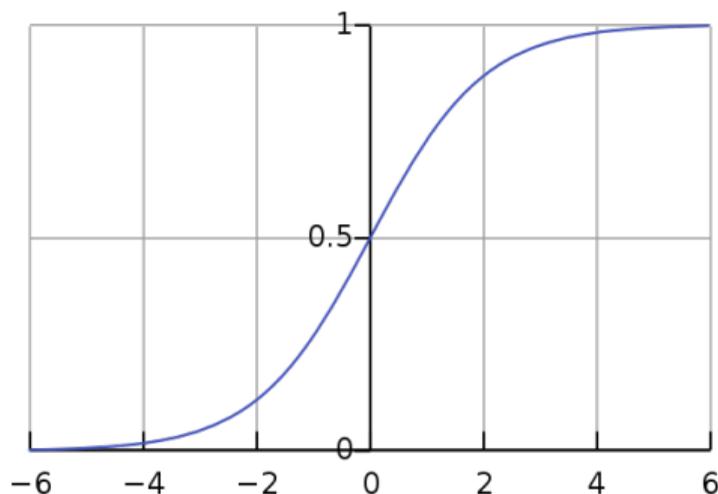
- Sigmoid function

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

- Input z is output of our regression

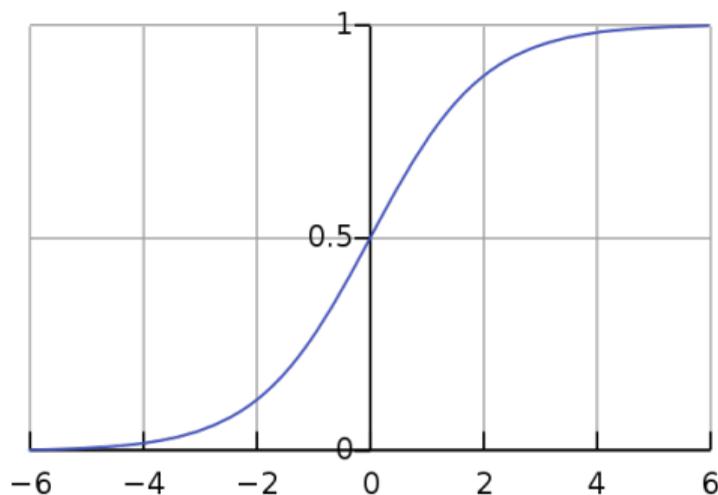
$$\sigma(z) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x_1 + \dots + \theta_k x_k)}}$$

- Adjust parameters to fix horizontal position and steepness of step



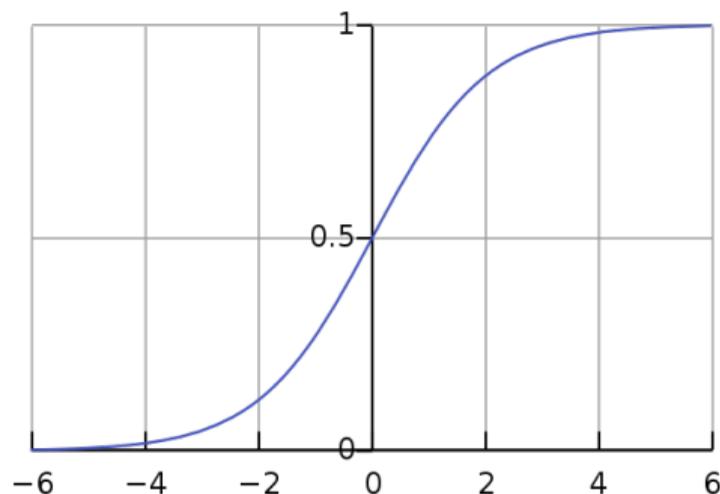
Logistic regression

- Compute the coefficients?
- Solve by gradient descent



Logistic regression

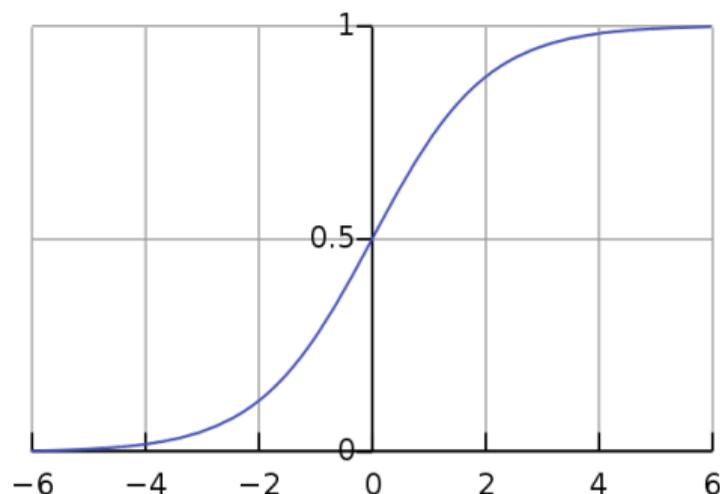
- Compute the coefficients?
- Solve by gradient descent
- Need derivatives to exist
 - Hence smooth sigmoid, not step function
 - $\sigma'(z) = \sigma(z)(1 - \sigma(z))$



Logistic regression

- Compute the coefficients?
- Solve by gradient descent
- Need derivatives to exist
 - Hence smooth sigmoid, not step function
 - $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Need a cost function to minimize

Max Likelihood



Loss function for logistic regression

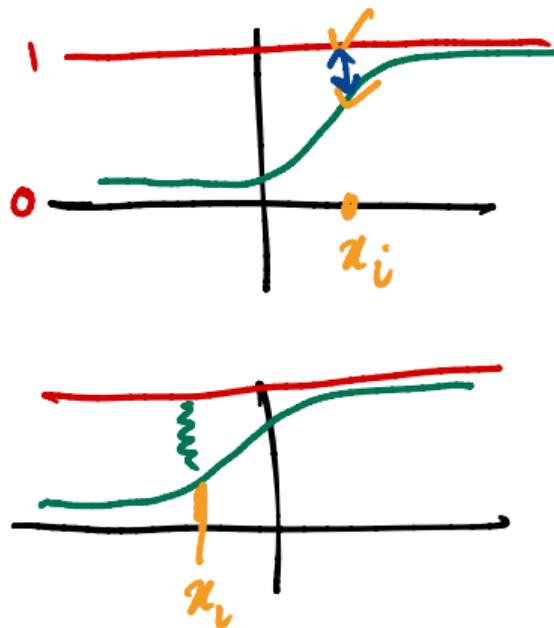
- Goal is to maximize log likelihood

Loss function for logistic regression

- Goal is to maximize log likelihood
- Let $h_{\theta}(x_i) = \sigma(z_i)$.

$$y_i = 1$$

$$y = 0$$

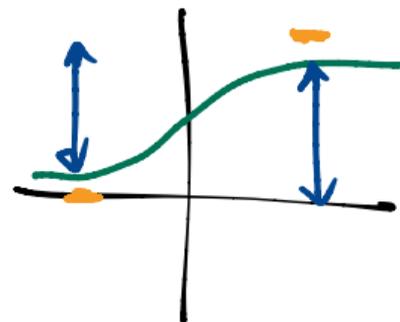


Loss function for logistic regression

- Goal is to maximize log likelihood
- Let $h_{\theta}(x_i) = \sigma(z_i)$. So, $P(y_i = 1 \mid x_i; \theta) = h_{\theta}(x_i)$,
 $P(y_i = 0 \mid x_i; \theta) = 1 - h_{\theta}(x_i)$
- Combine as $P(y_i \mid x_i; \theta) = \underbrace{h_{\theta}(x_i)^{y_i}}_{\text{if } y_i=1} \cdot \underbrace{(1 - h_{\theta}(x_i))^{1-y_i}}_{\text{if } y_i=0}$

Loss function for logistic regression

- Goal is to maximize log likelihood
- Let $h_{\theta}(x_i) = \sigma(z_i)$. So, $P(y_i = 1 | x_i; \theta) = h_{\theta}(x_i)$,
 $P(y_i = 0 | x_i; \theta) = 1 - h_{\theta}(x_i)$]
- Combine as $P(y_i | x_i; \theta) = h_{\theta}(x_i)^{y_i} \cdot (1 - h_{\theta}(x_i))^{1-y_i}$
- Likelihood: $\mathcal{L}(\theta) = \prod_{i=1}^n h_{\theta}(x_i)^{y_i} \cdot (1 - h_{\theta}(x_i))^{1-y_i}$



Loss function for logistic regression

- Goal is to maximize log likelihood

- Let $h_{\theta}(x_i) = \sigma(z_i)$. So, $P(y_i = 1 \mid x_i; \theta) = h_{\theta}(x_i)$,
 $P(y_i = 0 \mid x_i; \theta) = 1 - h_{\theta}(x_i)$

- Combine as $P(y_i \mid x_i; \theta) = h_{\theta}(x_i)^{y_i} \cdot (1 - h_{\theta}(x_i))^{1-y_i}$

- Likelihood: $\mathcal{L}(\theta) = \prod_{i=1}^n h_{\theta}(x_i)^{y_i} \cdot (1 - h_{\theta}(x_i))^{1-y_i}$

- Log-likelihood: $\ell(\theta) = \sum_{i=1}^n y_i \log h_{\theta}(x_i) + (1 - y_i) \log(1 - h_{\theta}(x_i))$

$$-(p_0 \log p_0 + p_i \log p_i)$$

Loss function for logistic regression

- Goal is to maximize log likelihood

- Let $h_\theta(x_i) = \sigma(z_i)$. So, $P(y_i = 1 \mid x_i; \theta) = h_\theta(x_i)$,
 $P(y_i = 0 \mid x_i; \theta) = 1 - h_\theta(x_i)$

- Combine as $P(y_i \mid x_i; \theta) = h_\theta(x_i)^{y_i} \cdot (1 - h_\theta(x_i))^{1-y_i}$

- Likelihood: $\mathcal{L}(\theta) = \prod_{i=1}^n h_\theta(x_i)^{y_i} \cdot (1 - h_\theta(x_i))^{1-y_i}$

- Log-likelihood: $\ell(\theta) = \sum_{i=1}^n y_i \log h_\theta(x_i) + (1 - y_i) \log(1 - h_\theta(x_i))$

- Minimize **cross entropy**: $-\sum_{i=1}^n y_i \log h_\theta(x_i) + (1 - y_i) \log(1 - h_\theta(x_i))$

MSE for logistic regression and gradient descent

- Suppose we take mean sum-squared error as the loss function.
- Consider two inputs $x = (x_1, x_2)$

$$C = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(z_i))^2, \text{ where } z_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

MSE for logistic regression and gradient descent

- Suppose we take mean sum-squared error as the loss function.
- Consider two inputs $x = (x_1, x_2)$

$$C = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(z_i))^2, \text{ where } z_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

- For gradient descent, we compute $\frac{\partial C}{\partial \theta_1}$, $\frac{\partial C}{\partial \theta_2}$, $\frac{\partial C}{\partial \theta_0}$

MSE for logistic regression and gradient descent

- Suppose we take mean sum-squared error as the loss function.
- Consider two inputs $x = (x_1, x_2)$

$$C = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(z_i))^2, \text{ where } z_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

- For gradient descent, we compute $\frac{\partial C}{\partial \theta_1}$, $\frac{\partial C}{\partial \theta_2}$, $\frac{\partial C}{\partial \theta_0}$

- For $j = 1, 2$,

$$\frac{\partial C}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (y_i - \sigma(z_i)) \cdot -\frac{\partial \sigma(z_i)}{\partial \theta_j}$$

MSE for logistic regression and gradient descent

- Suppose we take mean sum-squared error as the loss function.
- Consider two inputs $x = (x_1, x_2)$

$$C = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(z_i))^2, \text{ where } z_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

- For gradient descent, we compute $\frac{\partial C}{\partial \theta_1}$, $\frac{\partial C}{\partial \theta_2}$, $\frac{\partial C}{\partial \theta_0}$

- For $j = 1, 2$,

$$\frac{\partial C}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (y_i - \sigma(z_i)) \cdot \frac{\partial \sigma(z_i)}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \frac{\partial \sigma(z_i)}{\partial z_i} \frac{\partial z_i}{\partial \theta_j}$$

MSE for logistic regression and gradient descent

- Suppose we take mean sum-squared error as the loss function.
- Consider two inputs $x = (x_1, x_2)$

$$C = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(z_i))^2, \text{ where } z_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

- For gradient descent, we compute $\frac{\partial C}{\partial \theta_1}$, $\frac{\partial C}{\partial \theta_2}$, $\frac{\partial C}{\partial \theta_0}$

- For $j = 1, 2$,

$$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= \frac{2}{n} \sum_{i=1}^n (y_i - \sigma(z_i)) \cdot -\frac{\partial \sigma(z_i)}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \frac{\partial \sigma(z_i)}{\partial z_i} \frac{\partial z_i}{\partial \theta_j} \\ &= \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i) x_{ij} \end{aligned}$$

MSE for logistic regression and gradient descent

- Suppose we take mean sum-squared error as the loss function.
- Consider two inputs $x = (x_1, x_2)$

$$C = \frac{1}{n} \sum_{i=1}^n (y_i - \sigma(z_i))^2, \text{ where } z_i = \theta_0 + \theta_1 x_{i1} + \theta_2 x_{i2}$$

- For gradient descent, we compute $\frac{\partial C}{\partial \theta_1}$, $\frac{\partial C}{\partial \theta_2}$, $\frac{\partial C}{\partial \theta_0}$

- For $j = 1, 2$,

$$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= \frac{2}{n} \sum_{i=1}^n (y_i - \sigma(z_i)) \cdot -\frac{\partial \sigma(z_i)}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \frac{\partial \sigma(z_i)}{\partial z_i} \frac{\partial z_i}{\partial \theta_j} \\ &= \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i) x_{ij} \end{aligned}$$

- $\frac{\partial C}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \frac{\partial \sigma(z_i)}{\partial z_i} \frac{\partial z_i}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i)$

MSE for logistic regression and gradient descent . . .

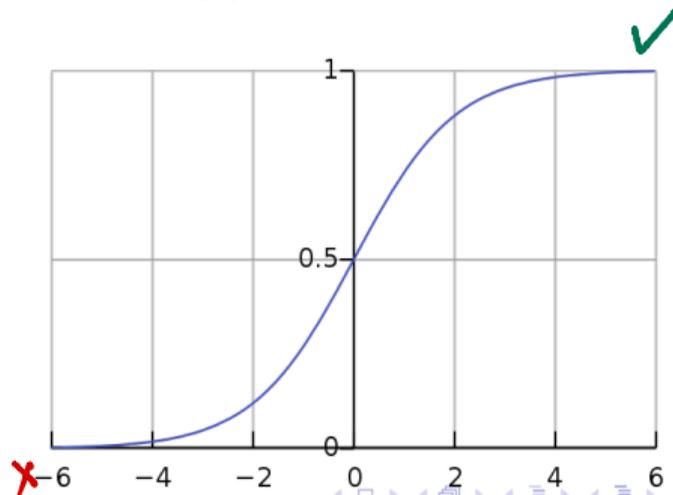
- For $j = 1, 2$, $\frac{\partial C}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i) x_j^i$, and $\frac{\partial C}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i)$
- Each term in $\frac{\partial C}{\partial \theta_1}$, $\frac{\partial C}{\partial \theta_2}$, $\frac{\partial C}{\partial \theta_0}$ is proportional to $\sigma'(z_i)$

MSE for logistic regression and gradient descent . . .

- For $j = 1, 2$, $\frac{\partial C}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i) x_j^i$, and $\frac{\partial C}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i)$
- Each term in $\frac{\partial C}{\partial \theta_1}$, $\frac{\partial C}{\partial \theta_2}$, $\frac{\partial C}{\partial \theta_0}$ is proportional to $\sigma'(z_i)$
- Ideally, gradient descent should take large steps when $\sigma(z) - y$ is large

MSE for logistic regression and gradient descent ...

- For $j = 1, 2$, $\frac{\partial C}{\partial \theta_j} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i) x_j^i$, and $\frac{\partial C}{\partial \theta_0} = \frac{2}{n} \sum_{i=1}^n (\sigma(z_i) - y_i) \sigma'(z_i)$
- Each term in $\frac{\partial C}{\partial \theta_1}$, $\frac{\partial C}{\partial \theta_2}$, $\frac{\partial C}{\partial \theta_0}$ is proportional to $\sigma'(z_i)$
- Ideally, gradient descent should take large steps when $\sigma(z) - y$ is large
- $\sigma(z)$ is flat at both extremes
- If $\sigma(z)$ is completely wrong, $\sigma(z) \approx (1 - y)$, we still have $\sigma'(z) \approx 0$
- Learning is slow even when current model is far from optimal



Cross entropy and gradient descent

- $C = -[y \ln(\sigma(z)) + (1 - y) \ln(1 - \sigma(z))]$

Cross entropy and gradient descent

- $C = -[y \ln(\sigma(z)) + (1 - y) \ln(1 - \sigma(z))]$

- $\frac{\partial C}{\partial \theta_j} = \frac{\partial C}{\partial \sigma} \frac{\partial \sigma}{\partial \theta_j}$

Cross entropy and gradient descent

- $C = -[y \ln(\sigma(z)) + (1 - y) \ln(1 - \sigma(z))]$

- $\frac{\partial C}{\partial \theta_j} = \frac{\partial C}{\partial \sigma} \frac{\partial \sigma}{\partial \theta_j} = - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \frac{\partial \sigma}{\partial \theta_j}$

Cross entropy and gradient descent

- $C = -[y \ln(\sigma(z)) + (1 - y) \ln(1 - \sigma(z))]$
- $$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= \frac{\partial C}{\partial \sigma} \frac{\partial \sigma}{\partial \theta_j} = - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \frac{\partial \sigma}{\partial \theta_j} \\ &= - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \frac{\partial \sigma}{\partial z} \frac{\partial z}{\partial \theta_j} \end{aligned}$$

Cross entropy and gradient descent

- $C = -[y \ln(\sigma(z)) + (1 - y) \ln(1 - \sigma(z))]$
- $$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= \frac{\partial C}{\partial \sigma} \frac{\partial \sigma}{\partial \theta_j} = - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \frac{\partial \sigma}{\partial \theta_j} \\ &= - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \frac{\partial \sigma}{\partial z} \frac{\partial z}{\partial \theta_j} \\ &= - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \sigma'(z) x_j \end{aligned}$$

Cross entropy and gradient descent

- $C = -[y \ln(\sigma(z)) + (1 - y) \ln(1 - \sigma(z))]$

- $$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= \frac{\partial C}{\partial \sigma} \frac{\partial \sigma}{\partial \theta_j} = - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \frac{\partial \sigma}{\partial \theta_j} \\ &= - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \frac{\partial \sigma}{\partial z} \frac{\partial z}{\partial \theta_j} \\ &= - \left[\frac{y}{\sigma(z)} - \frac{1 - y}{1 - \sigma(z)} \right] \sigma'(z) x_j \\ &= - \left[\frac{y(1 - \sigma(z)) - (1 - y)\sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma'(z) x_j \end{aligned}$$

Cross entropy and gradient descent . . .

- $\frac{\partial C}{\partial \theta_j} = - \left[\frac{y(1 - \sigma(z)) - (1 - y)\sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma'(z)x_j$
- Recall that $\sigma'(z) = \sigma(z)(1 - \sigma(z))$

Cross entropy and gradient descent ...

- $\frac{\partial C}{\partial \theta_j} = - \left[\frac{y(1 - \sigma(z)) - (1 - y)\sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma'(z)x_j$
- Recall that $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Therefore, $\frac{\partial C}{\partial \theta_j} = -[y(1 - \sigma(z)) - (1 - y)\sigma(z)]x_j$

Cross entropy and gradient descent ...

- $\frac{\partial C}{\partial \theta_j} = - \left[\frac{y(1 - \sigma(z)) - (1 - y)\sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma'(z)x_j$
- Recall that $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Therefore, $\frac{\partial C}{\partial \theta_j} = -[y(1 - \sigma(z)) - (1 - y)\sigma(z)]x_j$
 $= -[y - y\sigma(z) - \sigma(z) + y\sigma(z)]x_j$

Cross entropy and gradient descent ...

- $\frac{\partial C}{\partial \theta_j} = - \left[\frac{y(1 - \sigma(z)) - (1 - y)\sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma'(z)x_j$
- Recall that $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Therefore,
$$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= -[y(1 - \sigma(z)) - (1 - y)\sigma(z)]x_j \\ &= -[y - y\sigma(z) - \sigma(z) + y\sigma(z)]x_j \\ &= (\sigma(z) - y)x_j \end{aligned}$$

Cross entropy and gradient descent ...

- $\frac{\partial C}{\partial \theta_j} = - \left[\frac{y(1 - \sigma(z)) - (1 - y)\sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma'(z)x_j$
- Recall that $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Therefore,
$$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= -[y(1 - \sigma(z)) - (1 - y)\sigma(z)]x_j \\ &= -[y - y\sigma(z) - \sigma(z) + y\sigma(z)]x_j \\ &= (\sigma(z) - y)x_j \end{aligned}$$
- Similarly, $\frac{\partial C}{\partial \theta_0} = (\sigma(z) - y)$

Cross entropy and gradient descent . . .

- $\frac{\partial C}{\partial \theta_j} = - \left[\frac{y(1 - \sigma(z)) - (1 - y)\sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma'(z)x_j$
- Recall that $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Therefore,
$$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= -[y(1 - \sigma(z)) - (1 - y)\sigma(z)]x_j \\ &= -[y - y\sigma(z) - \sigma(z) + y\sigma(z)]x_j \\ &= (\sigma(z) - y)x_j \end{aligned}$$
- Similarly, $\frac{\partial C}{\partial \theta_0} = (\sigma(z) - y)$
- Thus, as we wanted, the gradient is proportional to $\sigma(z) - y$

Cross entropy and gradient descent ...

- $\frac{\partial C}{\partial \theta_j} = - \left[\frac{y(1 - \sigma(z)) - (1 - y)\sigma(z)}{\sigma(z)(1 - \sigma(z))} \right] \sigma'(z)x_j$
- Recall that $\sigma'(z) = \sigma(z)(1 - \sigma(z))$
- Therefore,
$$\begin{aligned} \frac{\partial C}{\partial \theta_j} &= -[y(1 - \sigma(z)) - (1 - y)\sigma(z)]x_j \\ &= -[y - y\sigma(z) - \sigma(z) + y\sigma(z)]x_j \\ &= (\sigma(z) - y)x_j \end{aligned}$$
- Similarly, $\frac{\partial C}{\partial \theta_0} = (\sigma(z) - y)$
- Thus, as we wanted, the gradient is proportional to $\sigma(z) - y$
- The greater the error, the faster the learning rate