CS 453X: Class 4

Jacob Whitehill

$$f_{\text{MSE}}(w_1, w_2) = \frac{1}{2}(x_1w_1 + x_2w_2 - y)^2$$

$$f_{\text{MSE}}(w_1, w_2) = \frac{1}{2} (x_1 w_1 + x_2 w_2 - y)^2$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_1}(w_1, w_2) = x_1 (x_1 w_1 + x_2 w_2 - y)$$

$$f_{\text{MSE}}(w_1, w_2) = \frac{1}{2} (x_1 w_1 + x_2 w_2 - y)^2$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_1}(w_1, w_2) = x_1 (x_1 w_1 + x_2 w_2 - y)$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_2}(w_1, w_2) = x_2 (x_1 w_1 + x_2 w_2 - y)$$

$$f_{\text{MSE}}(w_1, w_2) = \frac{1}{2}(x_1w_1 + x_2w_2 - y)^2$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_1}(w_1, w_2) = x_1(x_1w_1 + x_2w_2 - y)$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_2}(w_1, w_2) = x_2(x_1w_1 + x_2w_2 - y)$$

$$\nabla_{\mathbf{w}} f_{\text{MSE}}(\mathbf{w}) = \begin{bmatrix} \frac{\partial f_{\text{MSE}}}{\partial w_1} \\ \frac{\partial f_{\text{MSE}}}{\partial w_2} \end{bmatrix}$$

$$f_{\text{MSE}}(w_1, w_2) = \frac{1}{2} (x_1 w_1 + x_2 w_2 - y)^2$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_1} (w_1, w_2) = x_1 (x_1 w_1 + x_2 w_2 - y)$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_2} (w_1, w_2) = x_2 (x_1 w_1 + x_2 w_2 - y)$$

$$\nabla_{\mathbf{w}} f_{\text{MSE}}(\mathbf{w}) = \begin{bmatrix} \frac{\partial f_{\text{MSE}}}{\partial w_1} \\ \frac{\partial f_{\text{MSE}}}{\partial w_2} \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} (x_1 w_1 + x_2 w_2 - y)$$

$$f_{\text{MSE}}(w_1, w_2) = \frac{1}{2}(x_1w_1 + x_2w_2 - y)^2$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_1}(w_1, w_2) = x_1(x_1w_1 + x_2w_2 - y)$$

$$\frac{\partial f_{\text{MSE}}}{\partial w_2}(w_1, w_2) = x_2(x_1w_1 + x_2w_2 - y)$$

$$\nabla_{\mathbf{w}} f_{\text{MSE}}(\mathbf{w}) = \begin{bmatrix} \frac{\partial f_{\text{MSE}}}{\partial w_1} \\ \frac{\partial f_{\text{MSE}}}{\partial w_2} \end{bmatrix}$$

$$= \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} (x_1w_1 + x_2w_2 - y)$$

$$= \mathbf{x}(\mathbf{x}^{\top}\mathbf{w} - y)$$

Suppose we have n images, each with just 2 pixels.

$$\hat{\mathbf{y}} = \mathbf{X}^{\mathsf{T}} \mathbf{w}$$

Suppose we have n images, each with just 2 pixels.

$$\hat{\mathbf{y}} = \mathbf{X}^{\mathsf{T}} \mathbf{w}$$

$$= \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \dots & \mathbf{x}_{1}^{(n)} \\ \mathbf{x}_{2}^{(1)} & \dots & \mathbf{x}_{2}^{(n)} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} w_{1} \\ w_{2} \end{bmatrix}$$

This is the index of the *image*.

Suppose we have n images, each with just 2 pixels.

$$\hat{\mathbf{y}} = \mathbf{X}^{\mathsf{T}} \mathbf{w} \\
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \dots & \mathbf{x}_{1}^{(n)} \\ \mathbf{x}_{2}^{(1)} & \dots & \mathbf{x}_{2}^{(n)} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} w_{1} \\ w_{2} \end{bmatrix}$$

This is the index of the *pixel*.

Suppose we have n images, each with just 2 pixels.

$$\hat{\mathbf{y}} = \mathbf{X}^{\mathsf{T}} \mathbf{w}
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \dots & \mathbf{x}_{1}^{(n)} \\ \mathbf{x}_{2}^{(1)} & \dots & \mathbf{x}_{2}^{(n)} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} w_{1} \\ w_{2} \end{bmatrix}
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \mathbf{x}_{2}^{(1)} \\ \vdots & \vdots \\ \mathbf{x}_{1}^{(n)} & \mathbf{x}_{2}^{(n)} \end{bmatrix} \begin{bmatrix} w_{1} \\ w_{2} \end{bmatrix}$$

Suppose we have n images, each with just 2 pixels.

$$\hat{\mathbf{y}} = \mathbf{X}^{\mathsf{T}} \mathbf{w}
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \dots & \mathbf{x}_{1}^{(n)} \\ \mathbf{x}_{2}^{(1)} & \dots & \mathbf{x}_{2}^{(n)} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} w_{1} \\ w_{2} \end{bmatrix}
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \mathbf{x}_{2}^{(1)} \\ \vdots & \vdots \\ \mathbf{x}_{1}^{(n)} & \mathbf{x}_{2}^{(n)} \end{bmatrix} \begin{bmatrix} w_{1} \\ w_{2} \end{bmatrix}
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} w_{1} + \mathbf{x}_{2}^{(1)} w_{2} \\ \vdots \\ \mathbf{x}_{1}^{(n)} w_{1} + \mathbf{x}_{2}^{(n)} w_{2} \end{bmatrix}$$

Suppose we have *n* images, each with just 2 pixels.

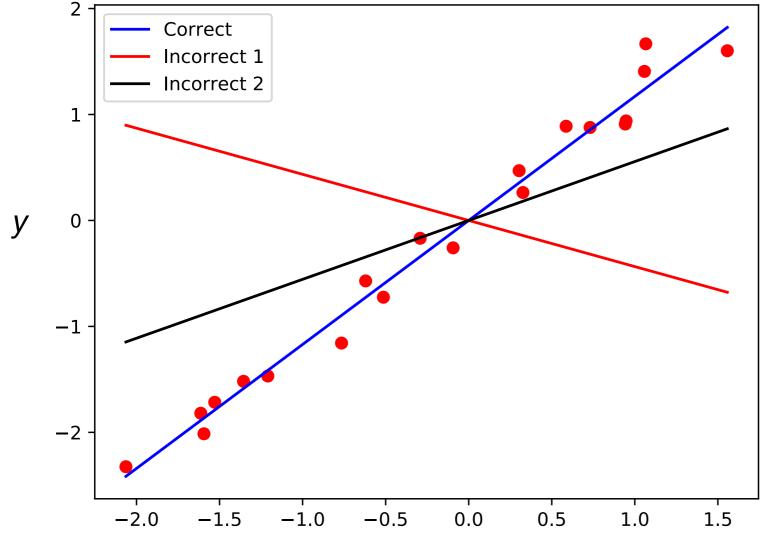
$$\hat{\mathbf{y}} = \mathbf{X}^{\mathsf{T}} \mathbf{w}
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \dots & \mathbf{x}_{1}^{(n)} \\ \mathbf{x}_{2}^{(1)} & \dots & \mathbf{x}_{2}^{(n)} \end{bmatrix}^{\mathsf{T}} \begin{bmatrix} w_{1} \\ w_{2} \end{bmatrix}
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} & \mathbf{x}_{2}^{(1)} \\ \vdots & \vdots \\ \mathbf{x}_{1}^{(n)} & \mathbf{x}_{2}^{(n)} \end{bmatrix} \begin{bmatrix} w_{1} \\ w_{2} \end{bmatrix}
= \begin{bmatrix} \mathbf{x}_{1}^{(1)} w_{1} + \mathbf{x}_{2}^{(1)} w_{2} \\ \vdots \\ \mathbf{x}_{1}^{(n)} w_{1} + \mathbf{x}_{2}^{(n)} w_{2} \end{bmatrix}$$

$$= \left[egin{array}{c} \hat{y}^{(1)} \ dots \ \hat{y}^{(n)} \end{array}
ight]$$

 $= \begin{bmatrix} \hat{y}^{(1)} \\ \vdots \\ \hat{y}^{(n)} \end{bmatrix}$ With just a single matrix-vector multiplication, we have computed our predictions for *all* of the *n* images simultaneously.

1-d example

• Linear regression finds the weight vector \mathbf{w} that minimizes the f_{MSE} . Here's an example where each \mathbf{x} is just 1-d...

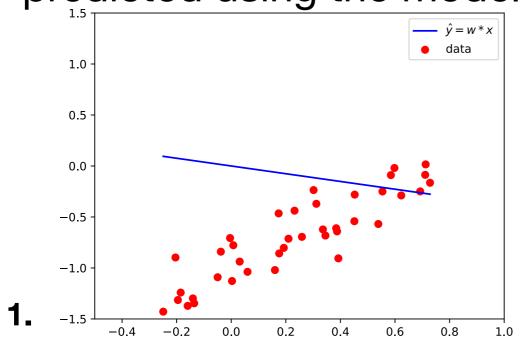


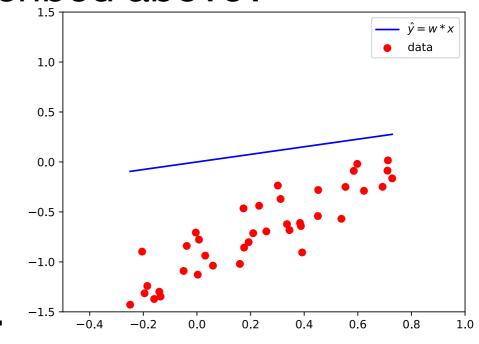
X

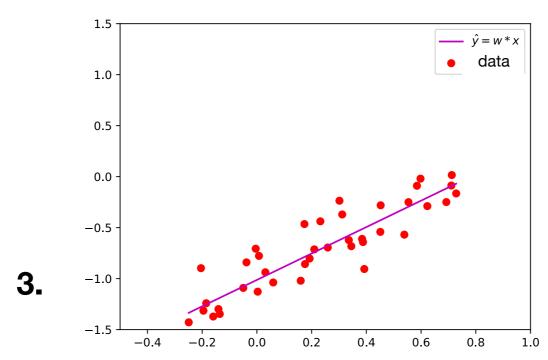
The best **w** is the one such that $f_{MSE}(\mathbf{y}, \, \hat{\mathbf{y}})$ is as small as possible, where each $\hat{y} = \mathbf{x}^{\mathsf{T}}\mathbf{w}$.

Exercise

 Which of the following regression lines would be predicted using the model described above?

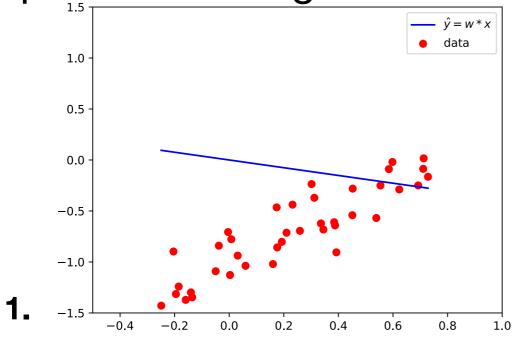






Exercise

 Which of the following regression lines would be predicted using the model described above?



Notice that the model enforces that (x,y)=(0,0) lie in the graph. Because of this constraint, the model learns the wrong slope to minimize the MSE.

 In order to account for target values with non-zero mean, we could add a bias term to our model:

$$\hat{y} = \mathbf{x}^{\top} \mathbf{w} + b$$

We could then compute the gradient w.r.t. both w and b and solve.

$$\nabla_{\mathbf{w}} f_{\text{MSE}}(\mathbf{y}, \hat{\mathbf{y}}; \mathbf{w}, b) = \nabla_{\mathbf{w}} \left[\frac{1}{2n} \sum_{i=1}^{n} \left(\mathbf{x}^{(i)^{\top}} \mathbf{w} + b - y^{(i)} \right)^{2} \right]$$

$$\nabla_{b} f_{\text{MSE}}(\mathbf{y}, \hat{\mathbf{y}}; \mathbf{w}, b) = \nabla_{b} \left[\frac{1}{2n} \sum_{i=1}^{n} \left(\mathbf{x}^{(i)^{\top}} \mathbf{w} + b - y^{(i)} \right)^{2} \right]$$

 Alternatively, we can implicitly include a bias term by augmenting each input vector x with a 1 at the end:

$$\tilde{\mathbf{x}} = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$$

 Correspondingly, our weight vector w will have an extra component (bias term) at the end.

$$\tilde{\mathbf{w}} = \left[\begin{array}{c} \mathbf{w} \\ b \end{array} \right]$$

To see why, notice that:

$$\hat{y} = \tilde{\mathbf{x}}^{\top} \tilde{\mathbf{w}}$$

$$= \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}^{\top} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$$

$$= \begin{bmatrix} \mathbf{x}^{\top} & 1 \end{bmatrix} \begin{bmatrix} \mathbf{w} \\ b \end{bmatrix}$$

$$= \mathbf{x}^{\top} \mathbf{w} + b$$

- We can find the optimal w and b based on all the training data using matrix notation.
- First define an augmented design matrix:

$$\tilde{\mathbf{X}} = \begin{bmatrix} \mathbf{x}^{(1)} & \dots & \mathbf{x}^{(n)} \\ 1 & \dots & 1 \end{bmatrix}$$

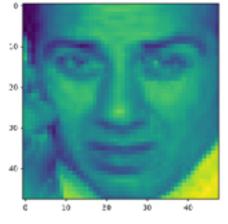
Then compute:

$$ilde{\mathbf{w}} = \left(ilde{\mathbf{X}} ilde{\mathbf{X}}^{ op}
ight)^{-1} ilde{\mathbf{X}} \mathbf{y}$$

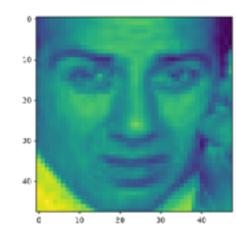
Example 1: age estimation

- Regress the age from 48x48 face images.
- Show demo...

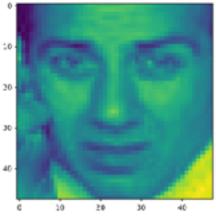
 How can we easily increase the number of training images?



Flip them left-right:



 How can we easily increase the number of training images?



- In numpy:
 - Let faces be a (n x 48 x 48) matrix containing n images.
 - Then facesFlipped = faces[:, :, ::-1] contains the left-right flipped images.

All *n* images.

- In numpy:
 - Let faces be a (n x 48 x 48) matrix containing n images.
 - Then facesFlipped = faces[:, :, ::-1] contains the left-right flipped images.

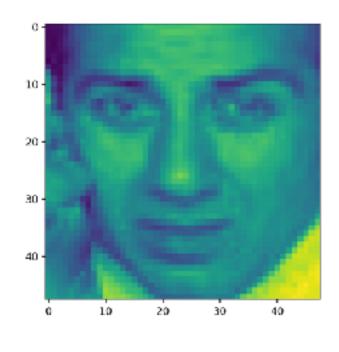
- In numpy:
 - Let faces be a (n x 48 x 48) matrix containing n images.
 - Then facesFlipped = faces[:, :, ::-1] contains the left-right flipped images.

All 48 columns but in reverse order.

- Avoid leakage of facial identity information:
 - Make sure that no "flipped" image in the testing set has an "unflipped" image in the training set (and viceversa).
- Data leakage: information in the training set which divulges information about the test set and which can bias the accuracy estimates.

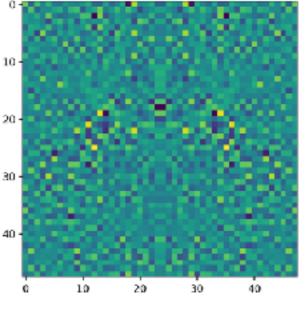
Learned weights

- Inspecting what the machine learned can be useful for debugging (and kinda fun to look at).
- For age estimation:



X





W

Higher temperatures associated with larger age values.