

NYU

Introduction to Robot Intelligence

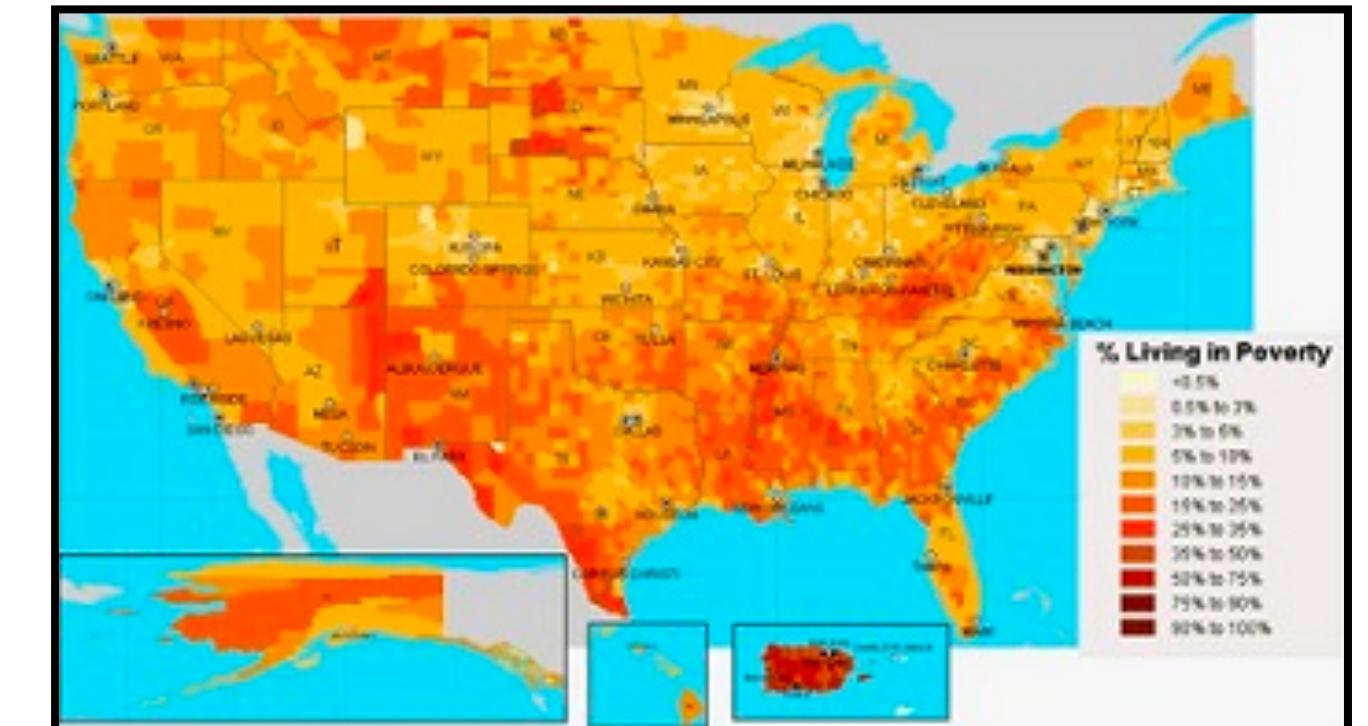
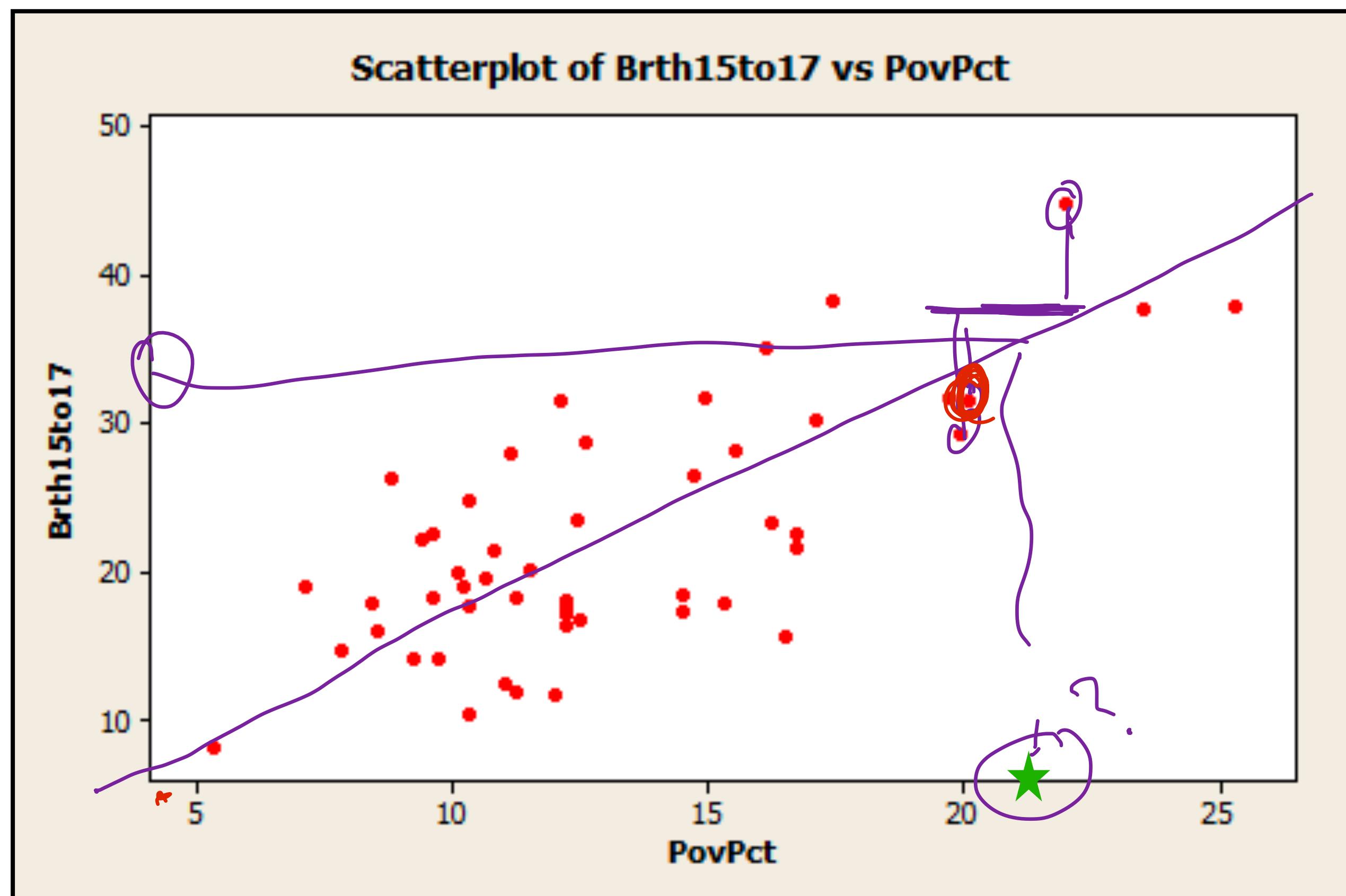
[Spring 2023]

Supervised Learning with SGD

February 7, 2023

Lerrel Pinto

Framing



<https://online.stat.psu.edu/stat462/node/101/>

Topics for today

- What is supervised learning?
- Introduction to Linear Regression
- Gradient Descent for supervised learning

Supervised Learning



Label: Cat



Label: Dog

Supervised Learning



Label: Cat

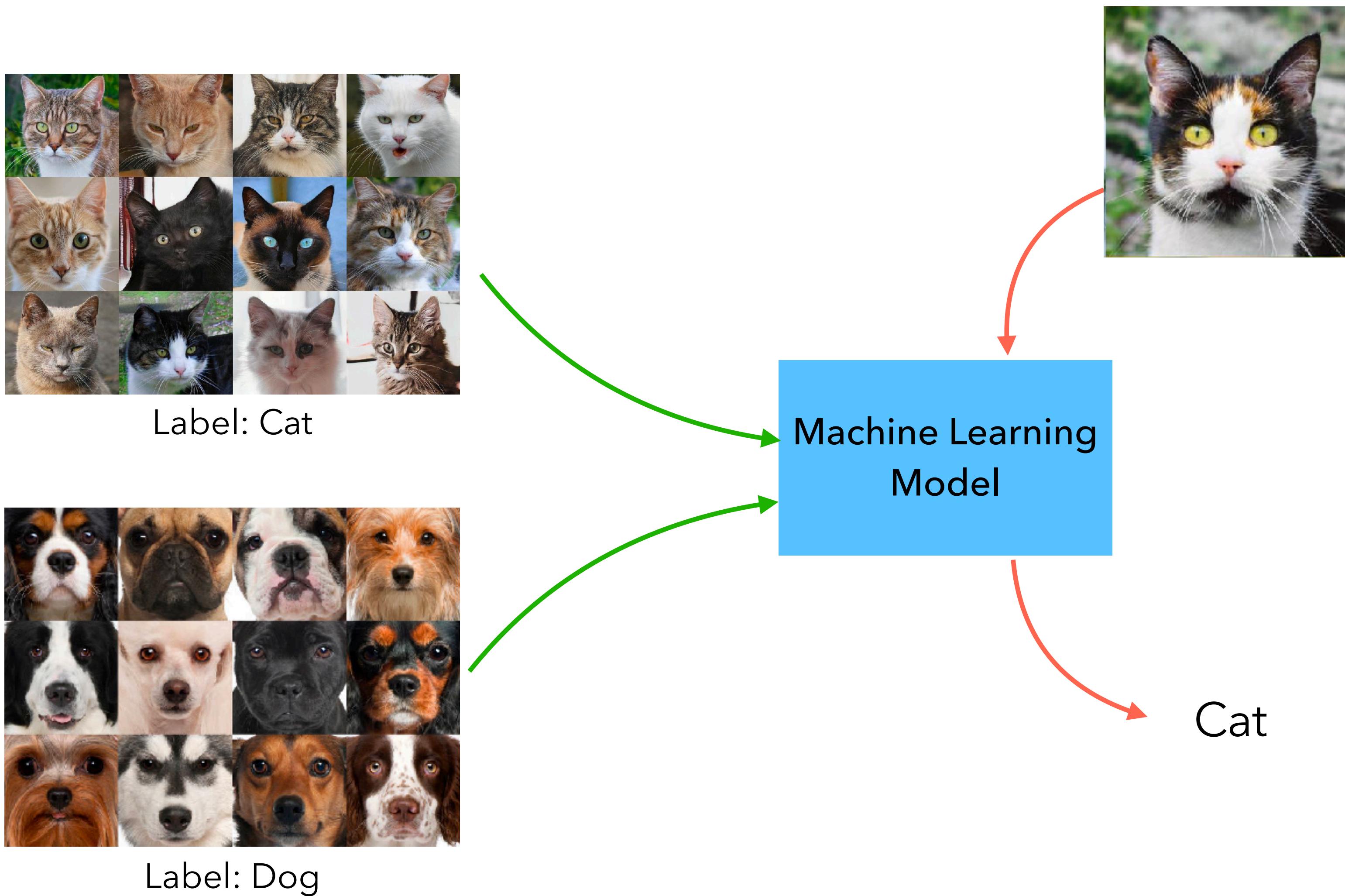


Label: Dog

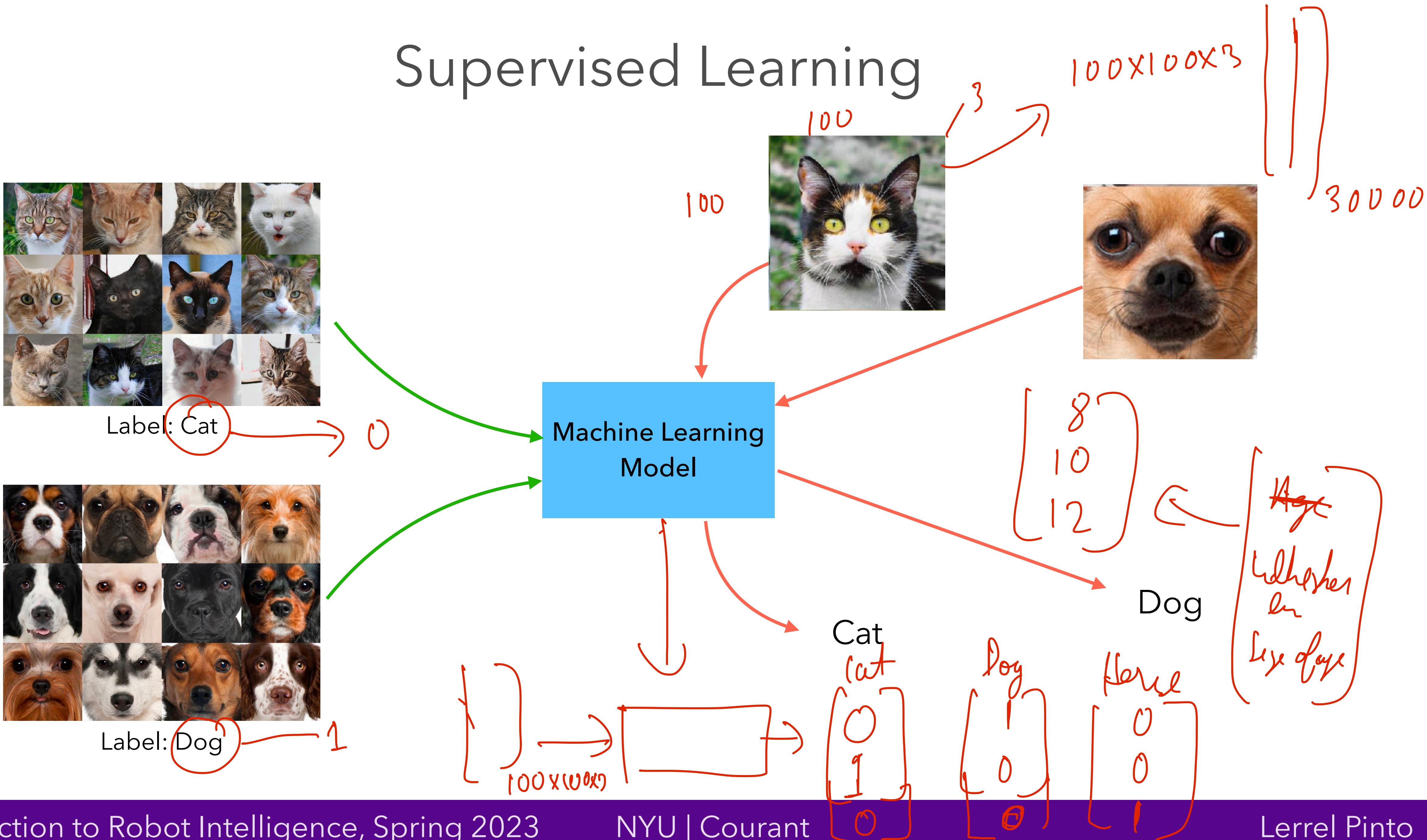
Machine Learning
Model



Supervised Learning



Supervised Learning



Examples of Supervised Learning

- Face detection. (Input: Image of face; Output: Name of person)
- Medical diagnosis. (Input: Test results; Output: Risk for heart disease)
- Video tagging. (Input: Uploaded video; Output: Caption for the video)
- Poverty Estimation. (Input: Night light; Output: Income level)
- Robotics. (Input: Sensor data; Output: Joint velocities)

Math Recap – Solving Linear Equations

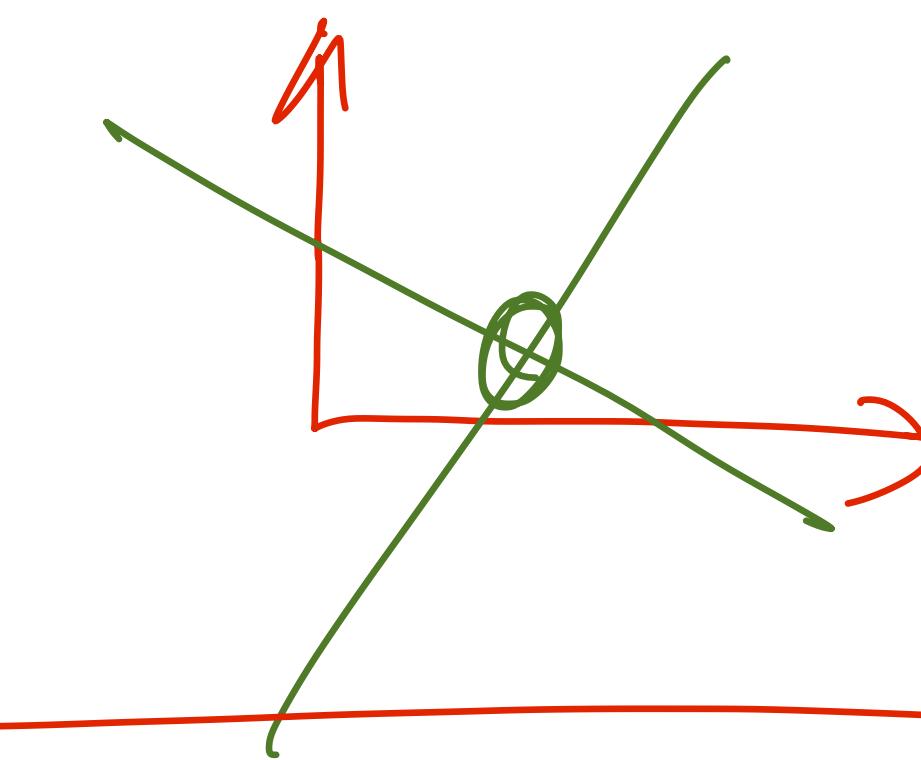
$$5 + \underline{8w_1} = 2$$

$$w_1 = \frac{2-5}{8} = -\frac{3}{8}$$

$$5(\textcircled{w}_1) + 8(\textcircled{w}_2) = 10$$

$$3w_1 + 9w_2 = 11$$

$$w_1 = \frac{10 - 8w_2}{5}$$



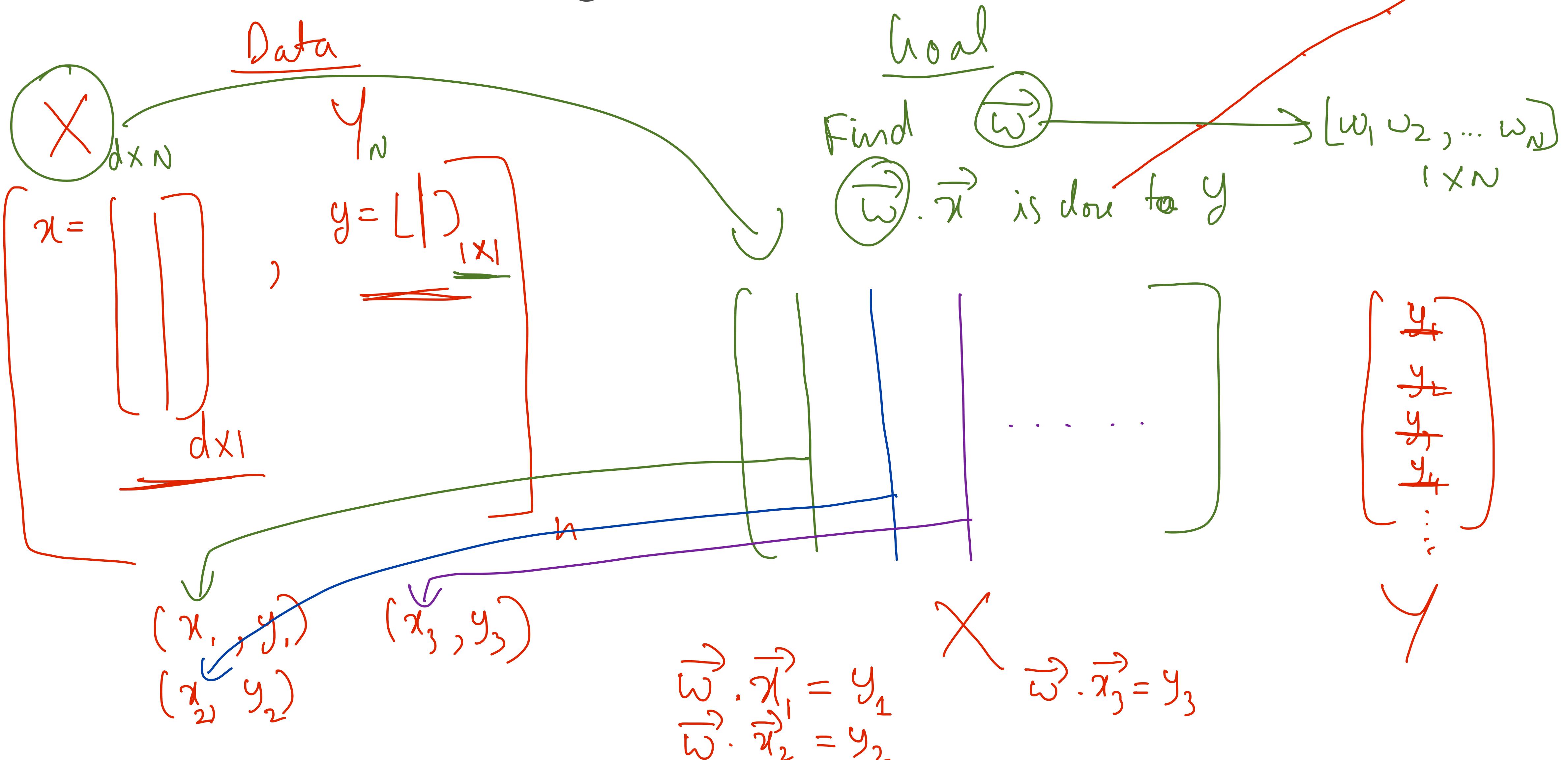
$$a_{11}w_1 + a_{12}w_2 + \dots + a_{1m}w_m = b_1$$

$$a_{21}w_1 + a_{22}w_2 + \dots + a_{2m}w_m = b_2$$

n equation
 m variable

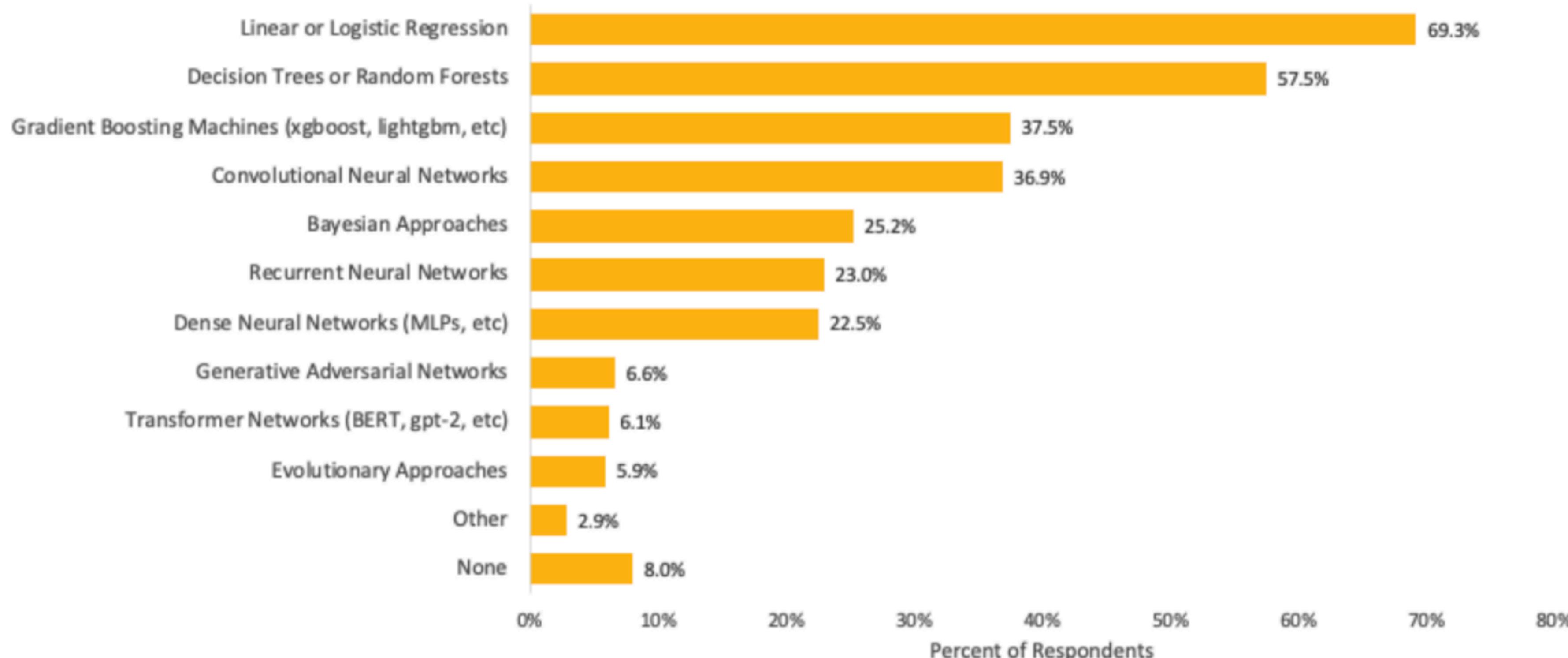
$$a_{n1}w_1 + a_{n2}w_2 + \dots + a_{nm}w_m = b_n$$

Linear Regression (formulation)



Linear Regression (why?)

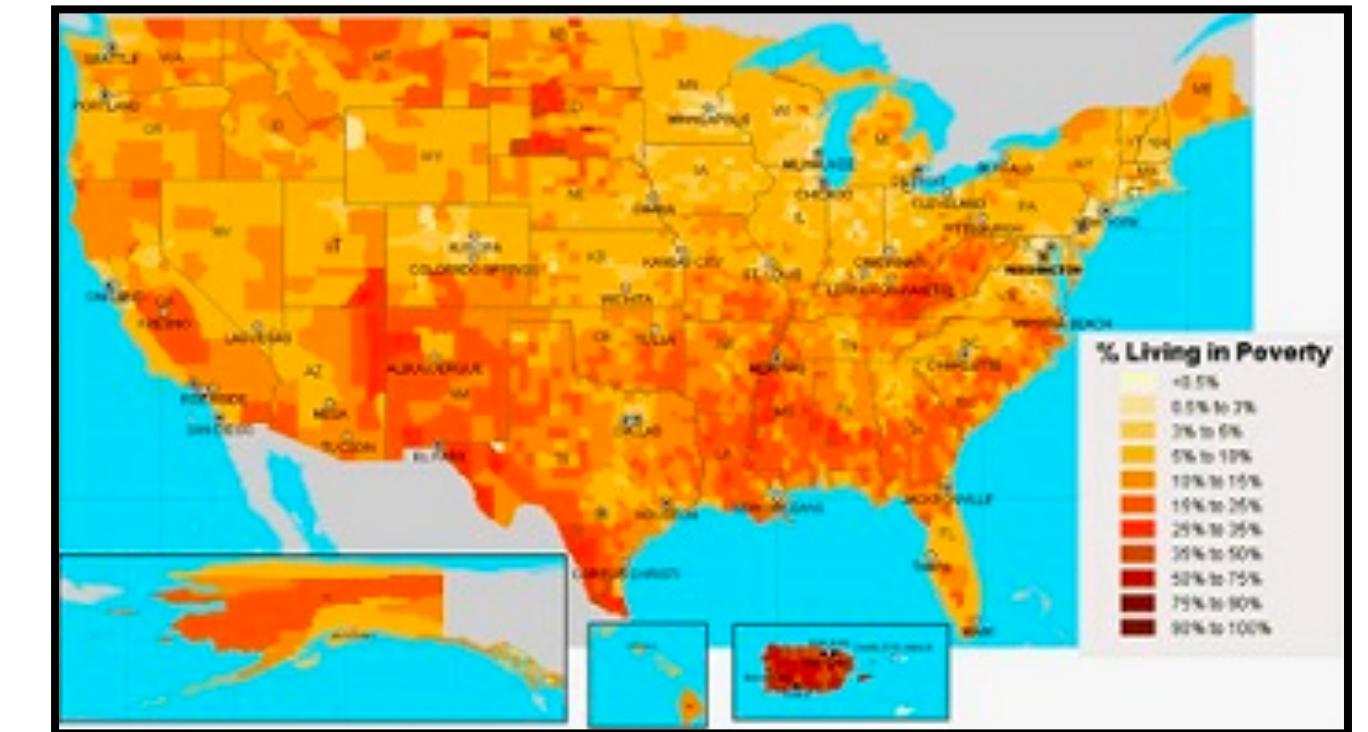
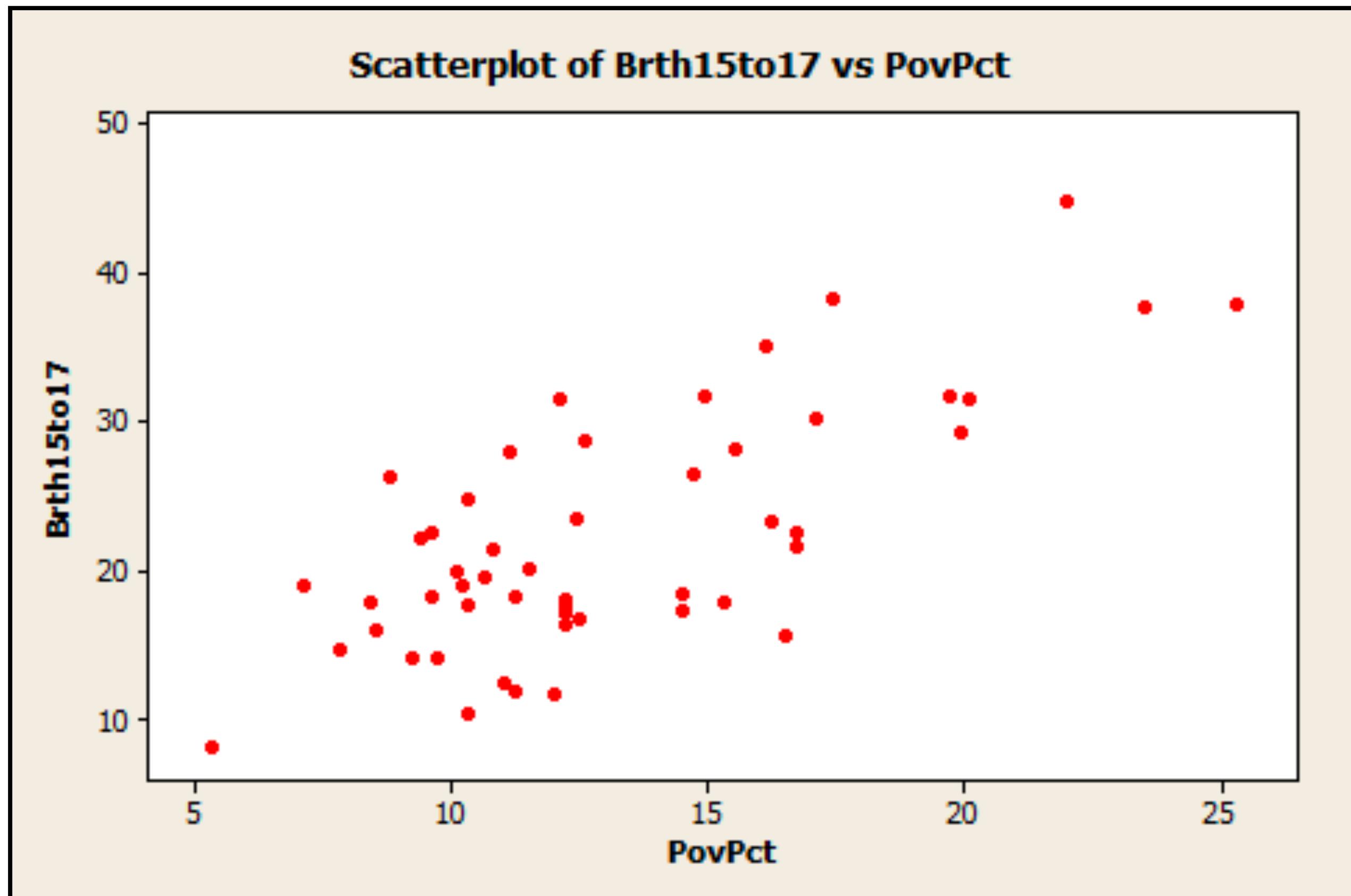
Which of the following ML algorithms do you use on a regular basis? (Select all that apply)



Note: Data are from the 2019 Kaggle ML and Data Science Survey. You can learn more about the study here: <https://www.kaggle.com/c/kaggle-survey-2019/data>.

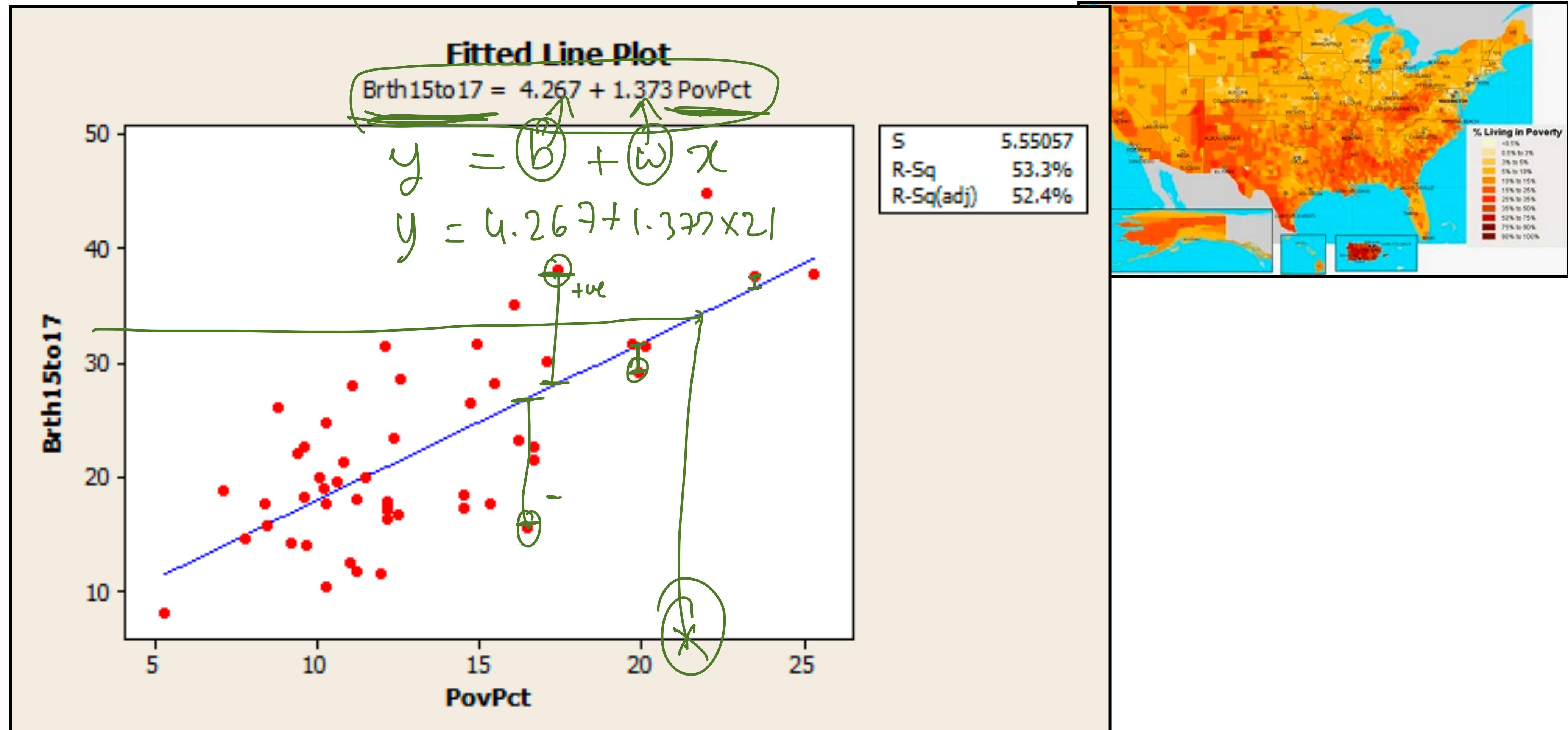
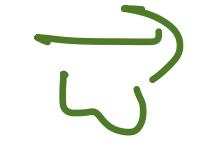
A total of 19717 respondents completed the survey; the percentages in the graph are based on a total of 14762 respondents who provided an answer to this question.

Linear Regression (example)



<https://online.stat.psu.edu/stat462/node/101/>

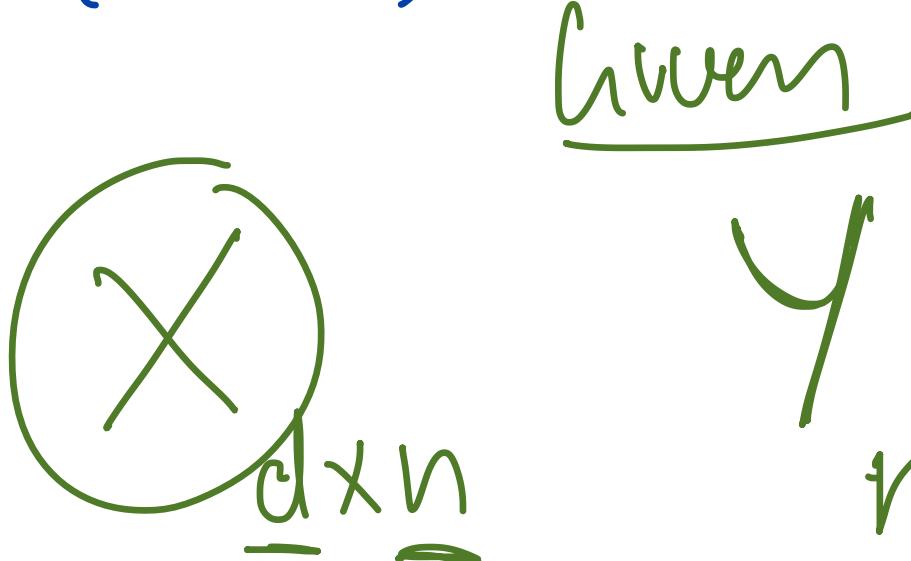
Linear Regression (example)



<https://online.stat.psu.edu/stat462/node/101/>

$$y = (ax+b)^2 \quad x = -\frac{b}{2a}$$

$$\frac{dy}{da} = 2(ax+b)a = 0$$



Linear Regression (objective)

$$\begin{aligned} & \text{minimize} \quad \sum_{j=1}^n e_j^2 \\ &= \sum_{i=1}^n (x_i^T w - y_i)^2 \\ &= \|x^T w - y\|^2 \end{aligned}$$

$$\left[\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \right] \left[\begin{array}{c} w_1 \\ w_2 \\ \vdots \\ w_n \end{array} \right] - \left[\begin{array}{c} y_1 \\ y_2 \\ \vdots \\ y_n \end{array} \right]$$

$$\begin{aligned} (e_1, e_2, \dots, e_n) \times \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} &= e_1^2 + e_2^2 + \dots + e_n^2 \\ x_1 \cdot w - y_1 \\ x_2 \cdot w - y_2 \\ \vdots \\ x_n \cdot w - y_n \end{pmatrix} &= \begin{pmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{pmatrix} \end{aligned}$$

Eg. 1

$$\begin{aligned} \vec{w} \cdot \vec{x}_1 \approx y_1 \rightarrow w^T x_1 \approx y_1 \\ \rightarrow x_1^T w \approx y_1 \end{aligned}$$

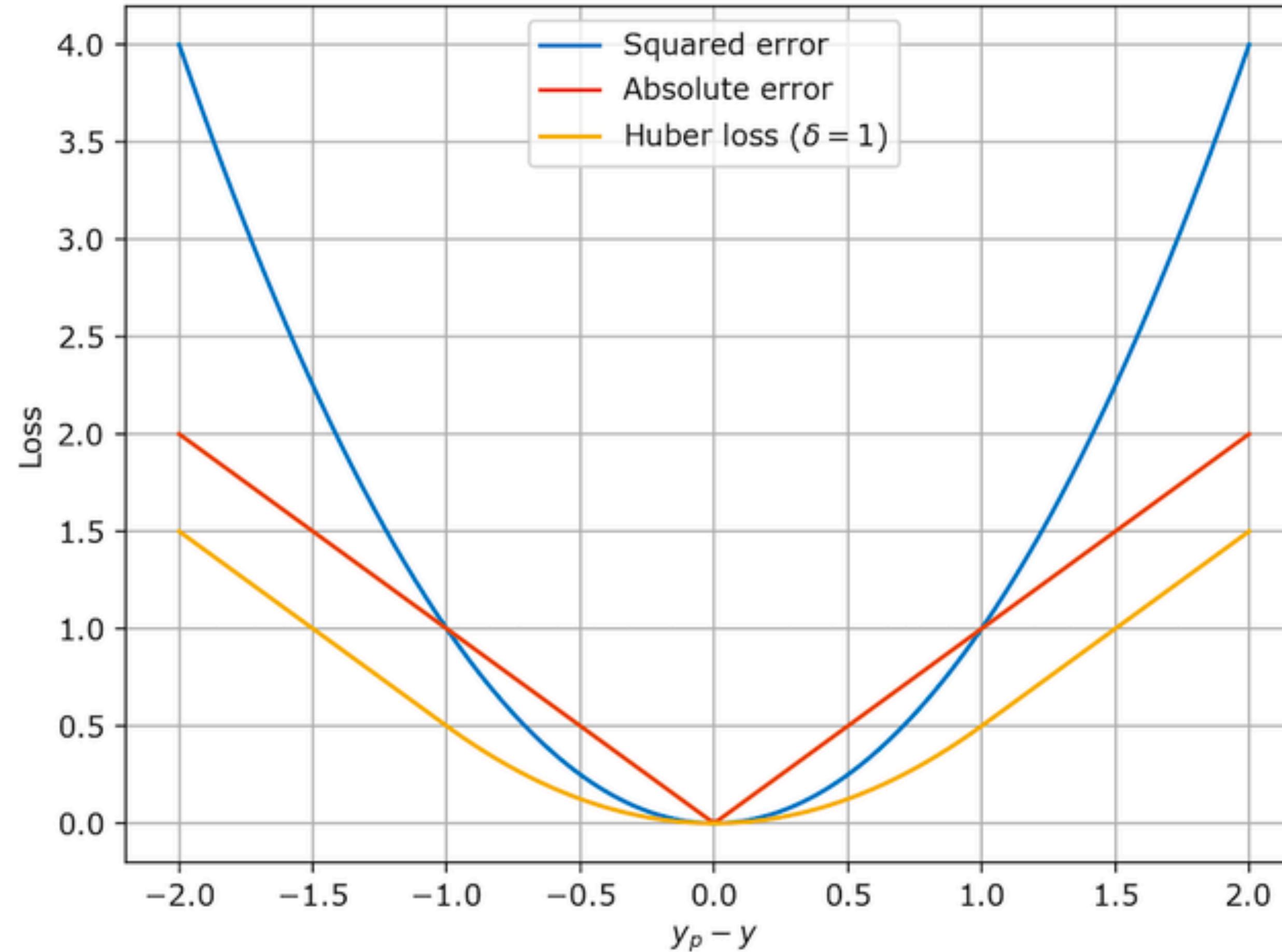
Eg. 2

$$\begin{aligned} e_1^2 &= (x_1^T w - y_1)^2 \\ e_2^2 &= (x_2^T w - y_2)^2 \end{aligned}$$

$$\begin{aligned} \vdots \\ \text{Eg. } n \quad e_n^2 &= (x_n^T w - y_n)^2 \\ \text{minimize} \quad \sum_{i=1}^n e_i^2 \end{aligned}$$

$$\begin{aligned} |e_1| &= e_1^2 \\ |e_2| &= e_2^2 \end{aligned}$$

Linear Regression (error function)



Linear Regression

- Input data: $X \in \mathbb{R}^{d \times n}$, $Y \in \mathbb{R}^n$, where $(\vec{x} \in \mathbb{R}^d, y \in \mathbb{R}^1)$ corresponds to a data point.

- $n \rightarrow$ # of data points, $d \rightarrow$ # of features / input dim.

- Goal: to find $\vec{w} \in \mathbb{R}^d$ such that $\langle \vec{w}, \vec{x} \rangle = y$

→ • Minimize $\|X^\top \vec{w} - Y\|^2$??

→ Solution: $\vec{w} = (XX^\top)^{-1}XY$

X^{-1}

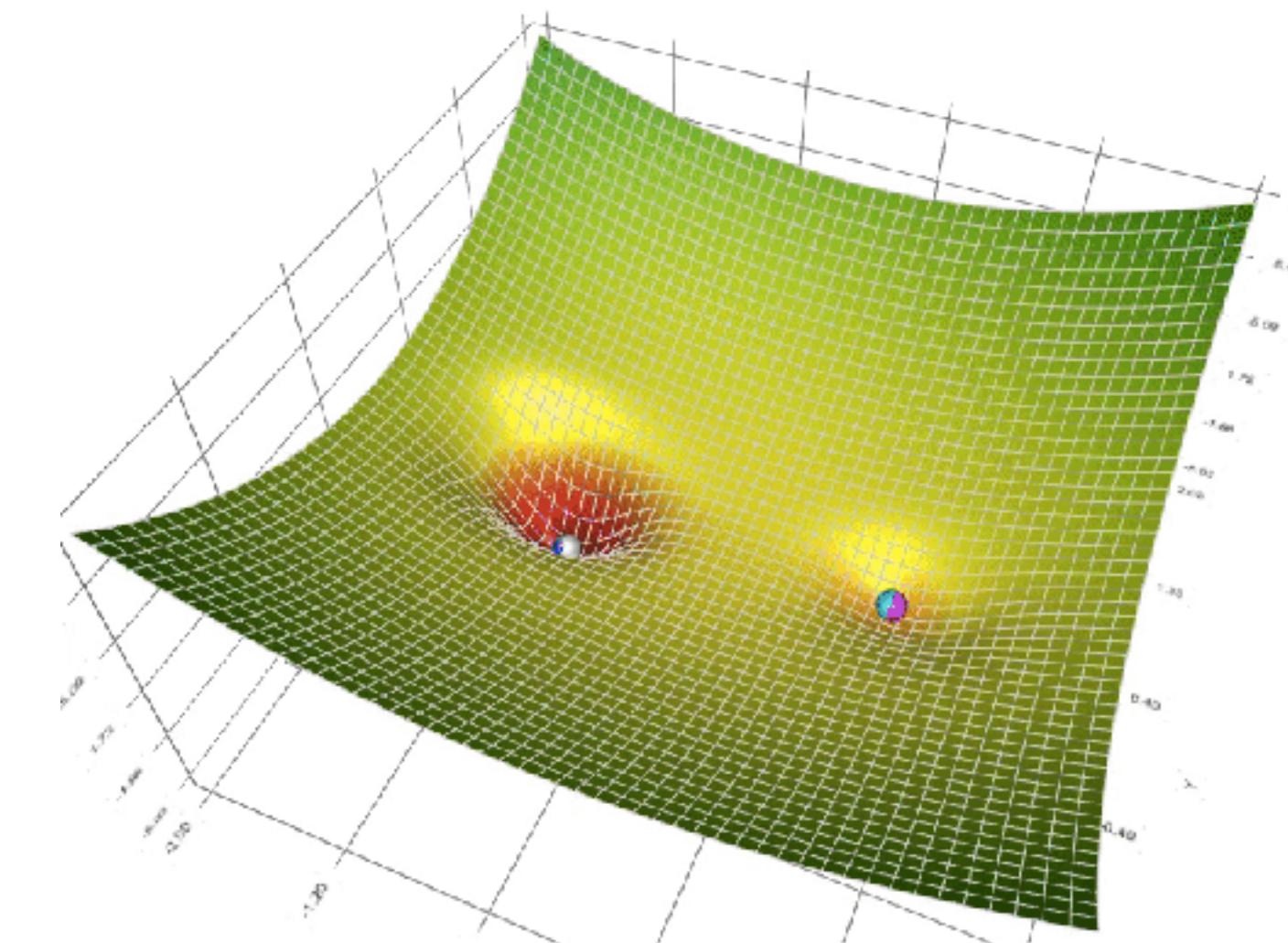
- Easy way to remember $\vec{w} = (X^\top)^+Y$

$$\begin{matrix} & d \times n & n \times d \\ R & R & R \\ & (d \times d) & \xrightarrow{\quad} \\ R & R & R \\ & d \times n & \end{matrix}$$

Recap: Gradient Descent

Gradient Descent Algorithm with multiple params

- Given: cost / loss/ objective function $f(\vec{\theta})$. Where $\vec{\theta} \in \mathbb{R}^d$.
- Goal: find $\vec{\theta}^*$ such that $f(\vec{\theta}^*) = \min_{\vec{\theta}} f(\vec{\theta})$.
- Gradient descent solution:
 - Start from initial guess $\vec{\theta}^0$ and learning rate α
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta})$
 - Repeat until change in θ is small, or maximum number of steps reached.



Linear Regression with Gradient Descent

Linear Regression with Gradient Descent

- Input data: $X \in \mathbb{R}^{d \times n}$, $Y \in \mathbb{R}^n$, where $(\vec{x} \in \mathbb{R}^d, y \in \mathbb{R}^1)$ corresponds to a data point.
 - $n \rightarrow$ # of data points, $d \rightarrow$ # of features / input dim.
- Goal: to find $\vec{w} \in \mathbb{R}^d$ such that $\langle \vec{w}, \vec{x} \rangle = y$
 - Minimize $\|X^\top \vec{w} - Y\|^2$
 - Analytic Solution: $\vec{w} = (XX^\top)^{-1}XY$

Linear Regression with Gradient Descent

- Input data: $X \in \mathbb{R}^{d \times n}$, $Y \in \mathbb{R}^n$, where $(\vec{x} \in \mathbb{R}^d, y \in \mathbb{R}^1)$ corresponds to a data point.
 - $n \rightarrow$ # of data points, $d \rightarrow$ # of features / input dim.
- Goal: to find $\vec{w} \in \mathbb{R}^d$ such that $\langle \vec{w}, \vec{x} \rangle = y$
 - Minimize $\|X^\top \vec{w} - Y\|^2$
 - Loss function $f(\vec{w}) = \|X^\top \vec{w} - Y\|^2 = (X^\top \vec{w} - Y)^\top (X^\top \vec{w} - Y)$

Linear Regression with Gradient Descent

- Loss function $f(\vec{w}) = \|X^T \vec{w} - Y\|^2 = (X^T \vec{w} - Y)^T (X^T \vec{w} - Y)$

Linear Regression with Gradient Descent

- Loss function $f(\vec{w}) = \|X^T \vec{w} - Y\|^2 = (X^T \vec{w} - Y)^T (X^T \vec{w} - Y)$

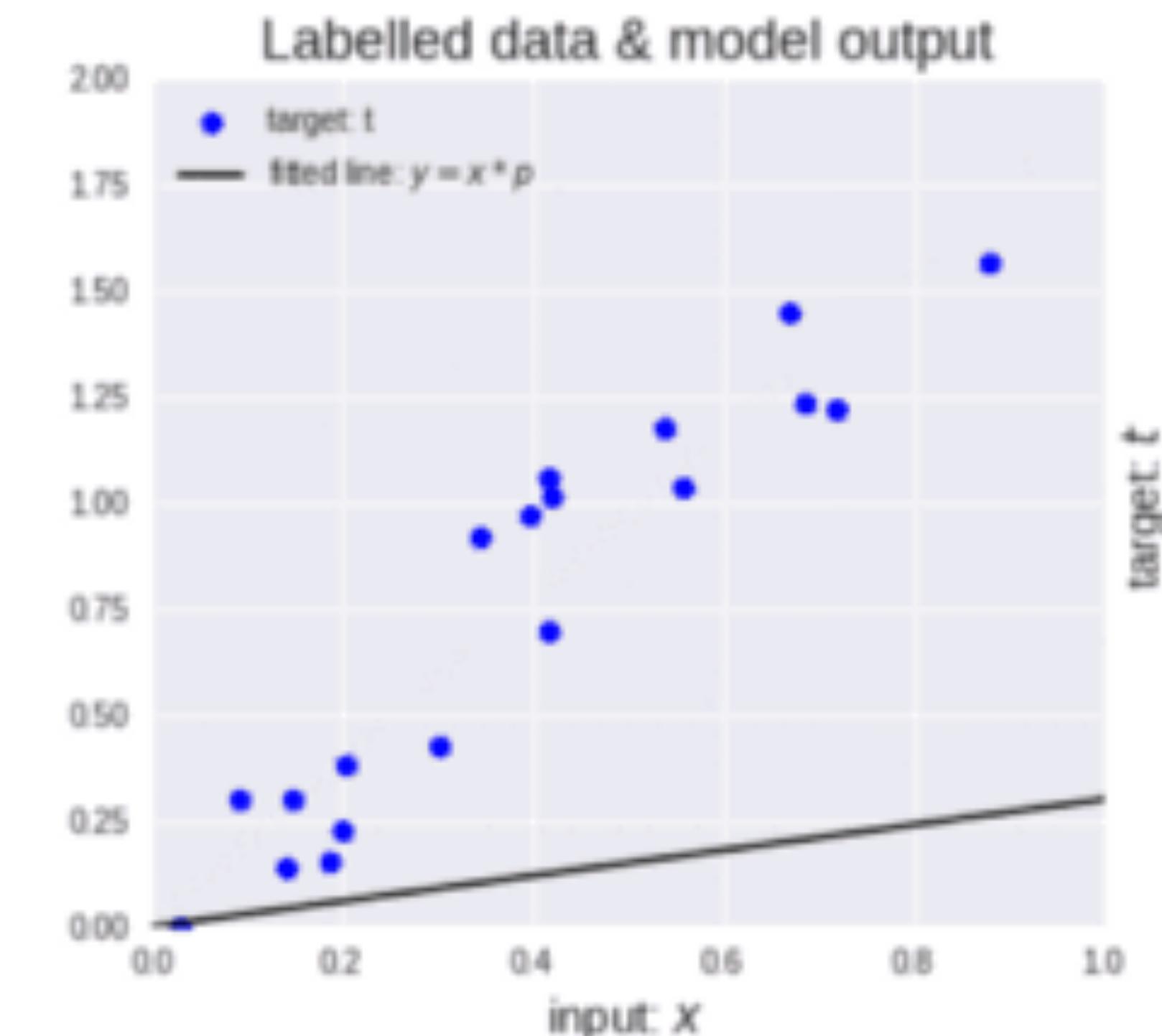
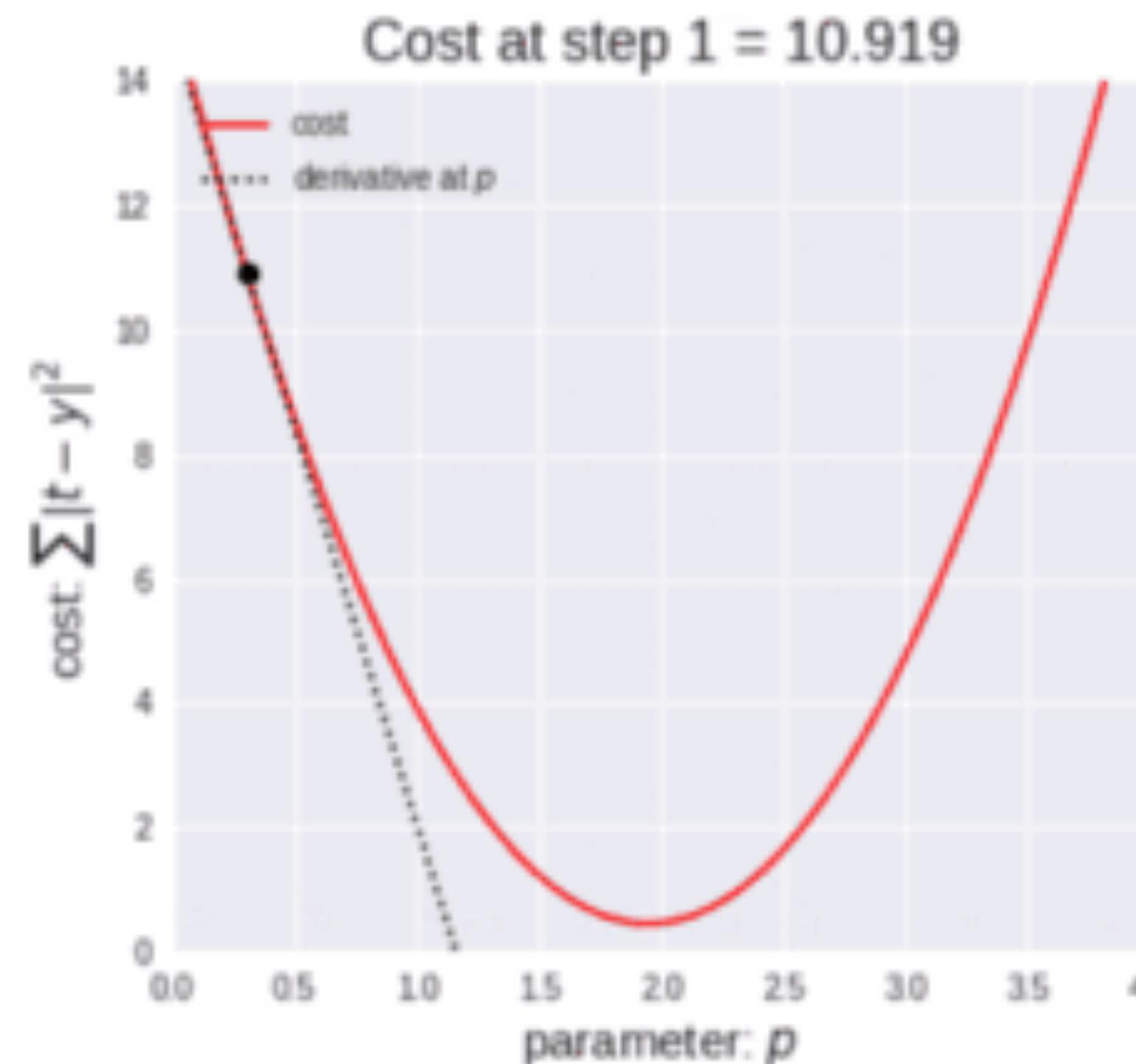
- Set $f(\vec{w}) = \sum_{i=1}^n (x^i \cdot \vec{w} - y^i)^2$

- Compute $\frac{\partial f}{\partial w_k} = \sum_{i=1}^n 2(\vec{w}^T x^i - y^i)x_k^i$

- Compute $\nabla f = [\frac{\partial f}{\partial w_1}, \dots, \frac{\partial f}{\partial w_k}, \dots, \frac{\partial f}{\partial w_d}]$

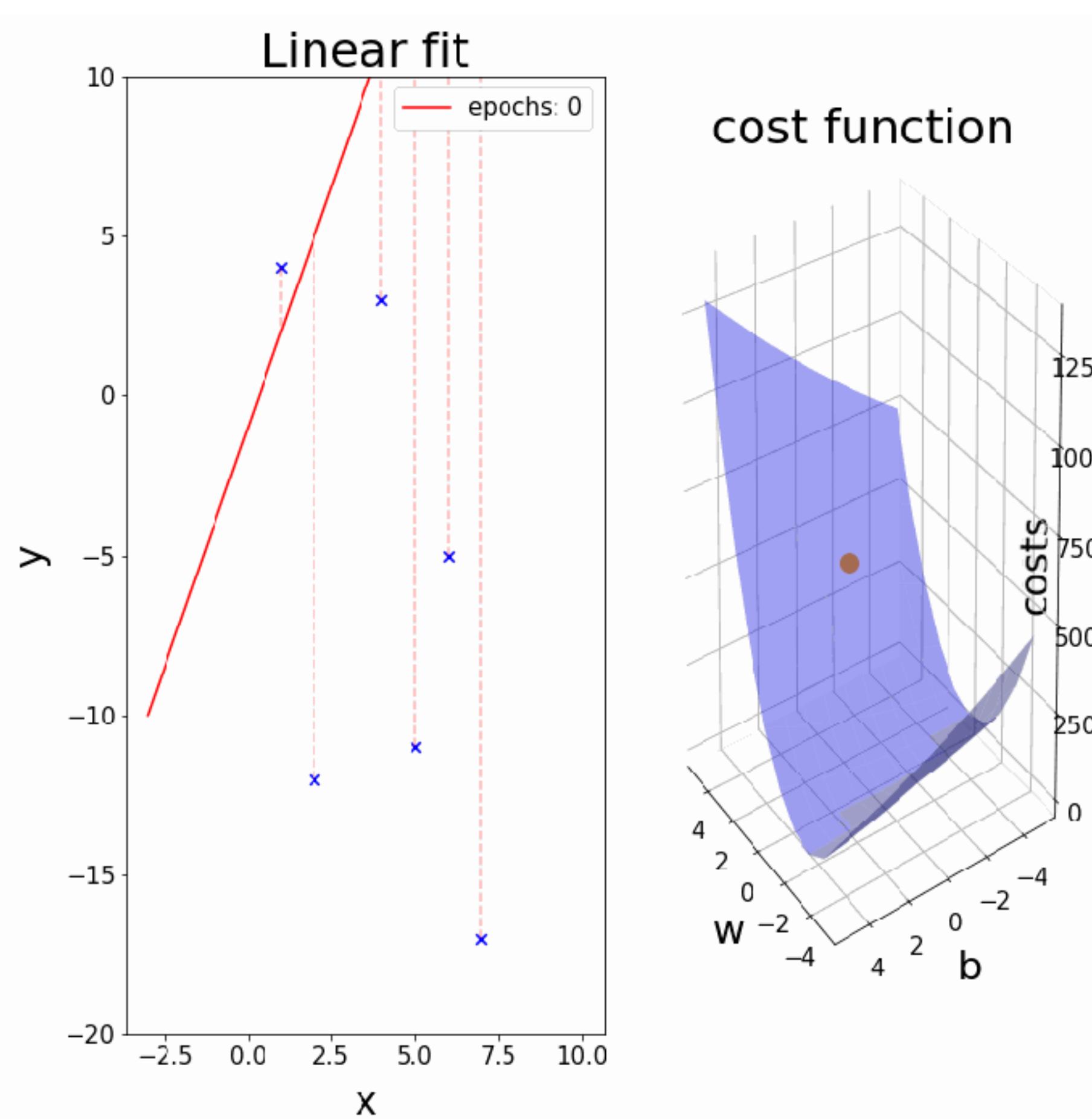
- Use gradient descent

Linear Regression with Gradient Descent



Credits: <https://www.kdnuggets.com/>

Linear Regression with Gradient Descent



Credits: Tobias Roeschl

Gradient Descent Algorithm for Machine Learning

Gradient Descent Algorithm for Machine Learning

- Given: cost / loss/ objective function $f(\vec{\theta}, \textcolor{blue}{D})$. Where $\vec{\theta} \in \mathbb{R}^d$.

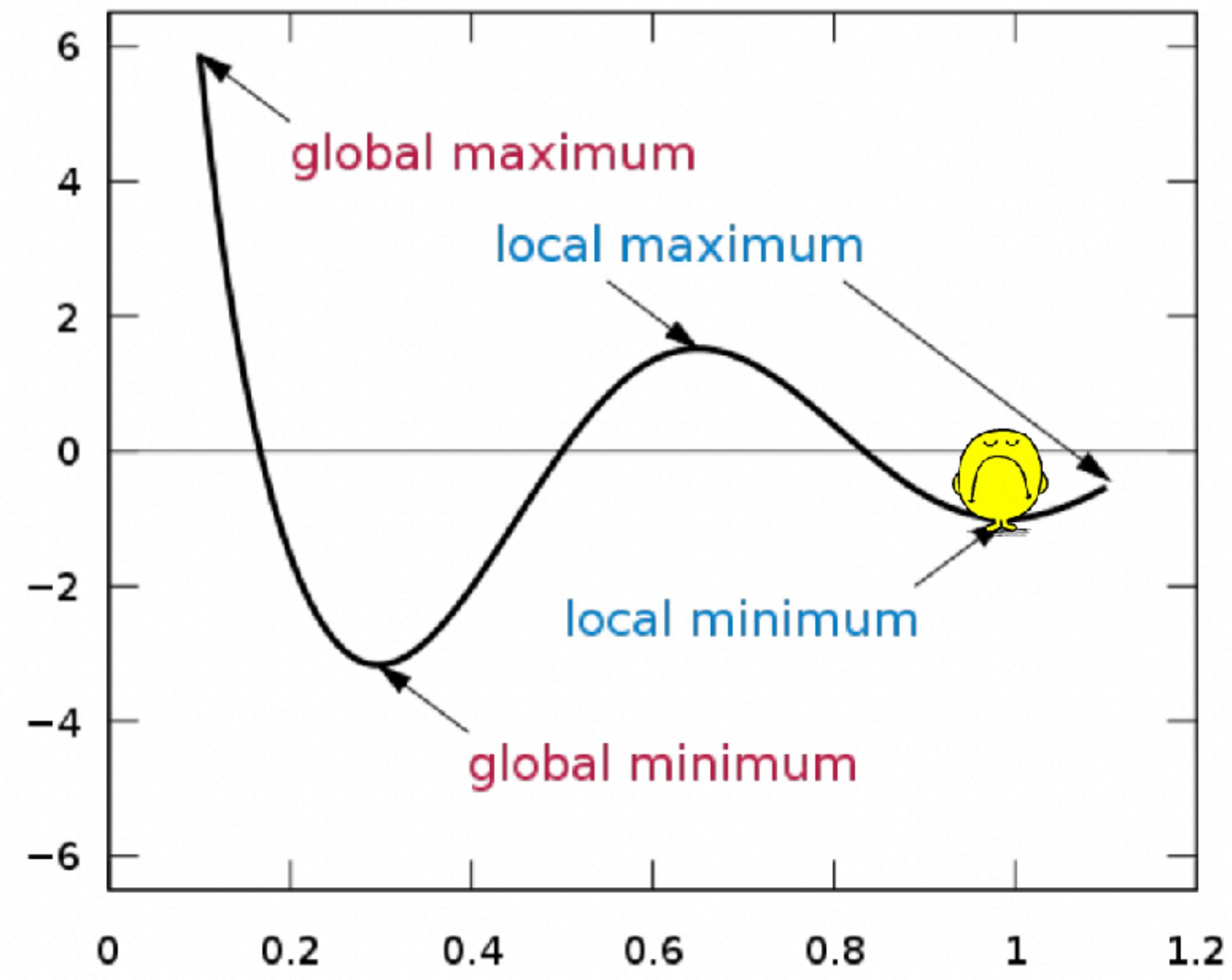
Gradient Descent Algorithm for Machine Learning

- Given: cost / loss/ objective function $f(\vec{\theta}, \textcolor{blue}{D})$. Where $\vec{\theta} \in \mathbb{R}^d$.
- Goal: find $\vec{\theta}^*$ such that $f(\vec{\theta}^*, \textcolor{blue}{D}) = \min_{\vec{\theta}} f(\vec{\theta}, \textcolor{blue}{D})$.

Gradient Descent Algorithm for Machine Learning

- Given: cost / loss/ objective function $f(\vec{\theta}, \textcolor{blue}{D})$. Where $\vec{\theta} \in \mathbb{R}^d$.
- Goal: find $\vec{\theta}^*$ such that $f(\vec{\theta}^*, \textcolor{blue}{D}) = \min_{\vec{\theta}} f(\vec{\theta}, \textcolor{blue}{D})$.
- Gradient descent solution:
 - Start from initial guess $\vec{\theta}^0$ and learning rate α
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta}, \textcolor{blue}{D})$
 - Repeat until change in θ is small, or maximum number of steps reached.

Key issue: Local minima



Credits: Michail Michailidis & Patrick Maiden

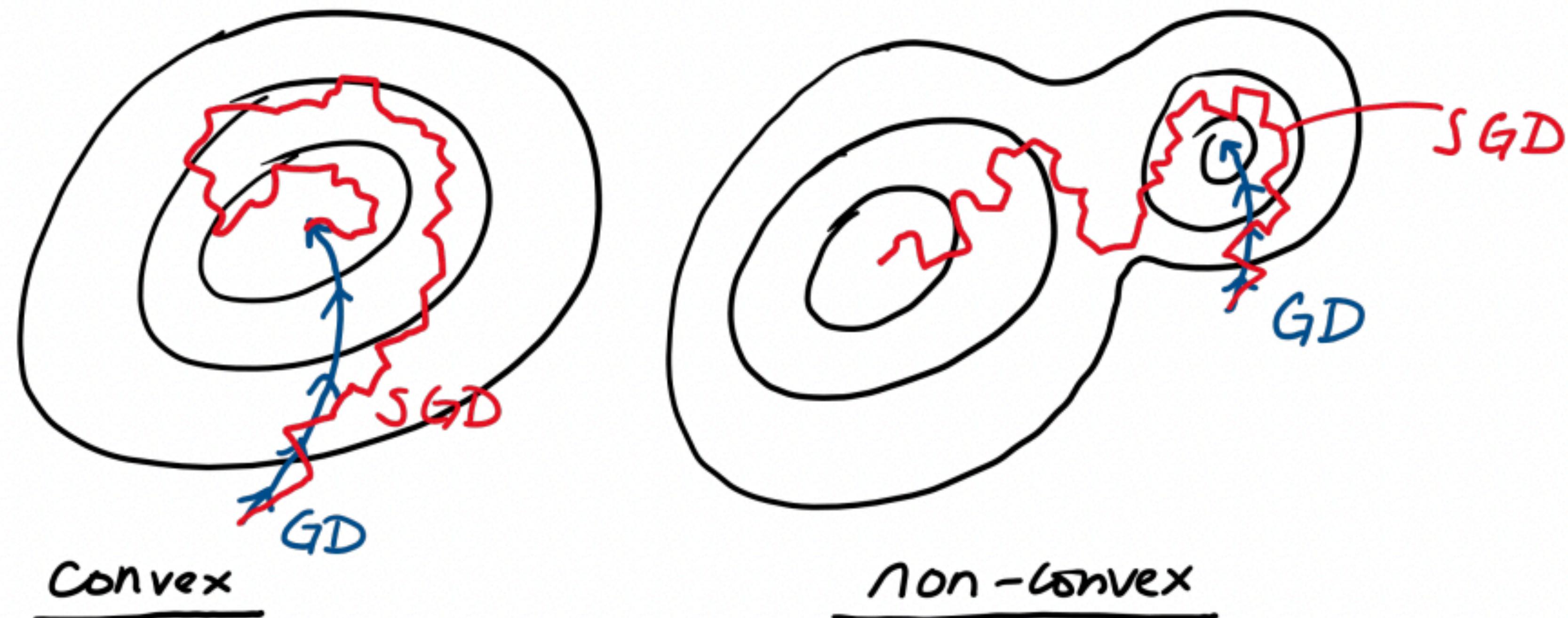
One solution: Stochastic Gradient Descent (SGD)

- Gradient descent solution:
 - Start from initial guess $\vec{\theta}^0$ and learning rate α
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta}, D)$

One solution: Stochastic Gradient Descent (SGD)

- Gradient descent solution:
 - Start from initial guess $\vec{\theta}^0$ and learning rate α
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta}, D)$
- SGD:
 - Sample single or multiple datapoints $d \sim D$
 - Update $\vec{\theta}^{i+1} \leftarrow \vec{\theta}^i - \alpha \nabla f(\vec{\theta}, d)$

One solution: Stochastic Gradient Descent (SGD)



Credits: Stanley Chan

Additional Reading

- <https://www.mit.edu/~6.s085/notes/lecture3.pdf>
- https://www.coconino.edu/resources/files/pdfs/academics/sabbatical-reports/kate-kozak/chapter_10.pdf
- Interactive tutorial: <https://uclaacm.github.io/gradient-descent-visualiser/>
- Book chapter: <https://www.cs.utah.edu/~jeffp/IDABook/T6-GD.pdf>
- SGD + variants: <https://ruder.io/optimizing-gradient-descent>

Questions?