

Introduction to Regressi

Srinandan (“Sri”) Dasmahapatra

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

COMP3223

Supervised Learning: labelled data

Compare labels with predictions

$d(\hat{y}_n, y_n)$: How far is prediction \hat{y}_n from actual data y_n ?

- Given data \mathcal{D} construct model $f(\cdot; \mathbf{w})$ such that the “**distance**” between model output and real “output” is small

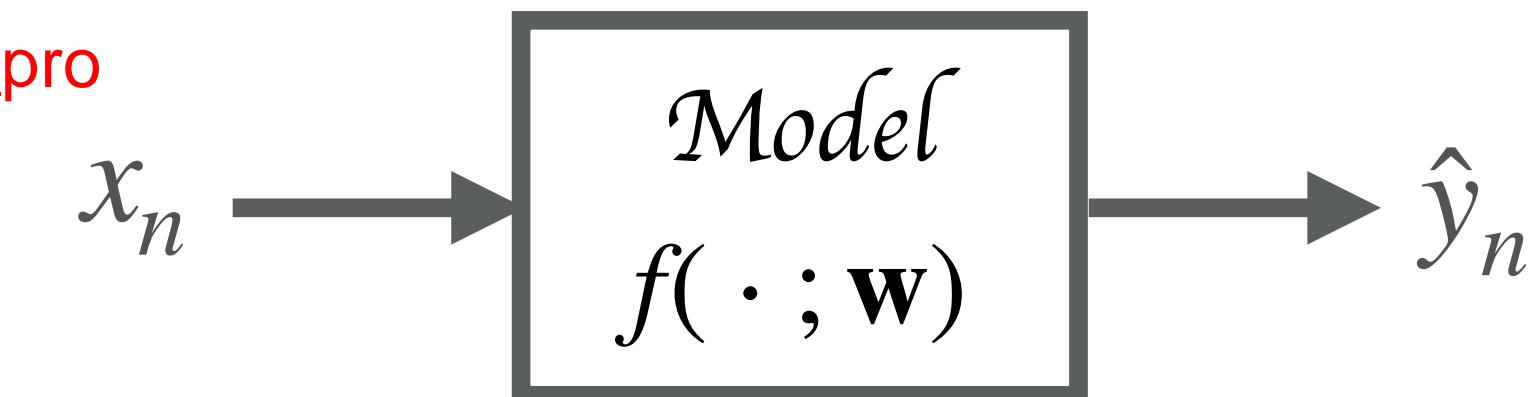
Assignment Project Exam Help
<https://eduassistpro.github.io/>

- Learning = model construction by

minimising loss $L(\mathbf{w}) = \sum_{n=1}^N d(\hat{y}_n, y_n)$

- $y = f(\cdot; \mathbf{w})$ continuous (e.g., **y=23.4**),
 $f(\cdot; \mathbf{w})$ is a **regression** model

$$\mathcal{D} := \{(x_n, y_n)\}, n = 1, \dots, N$$



$$\hat{y}_n = f(x_n; \mathbf{w})$$



Core idea in ML: Reduce mismatch between model prediction and data

Squared residual loss

- Regression models minimise **residuals** – deviations of model predictions from outputs in training data

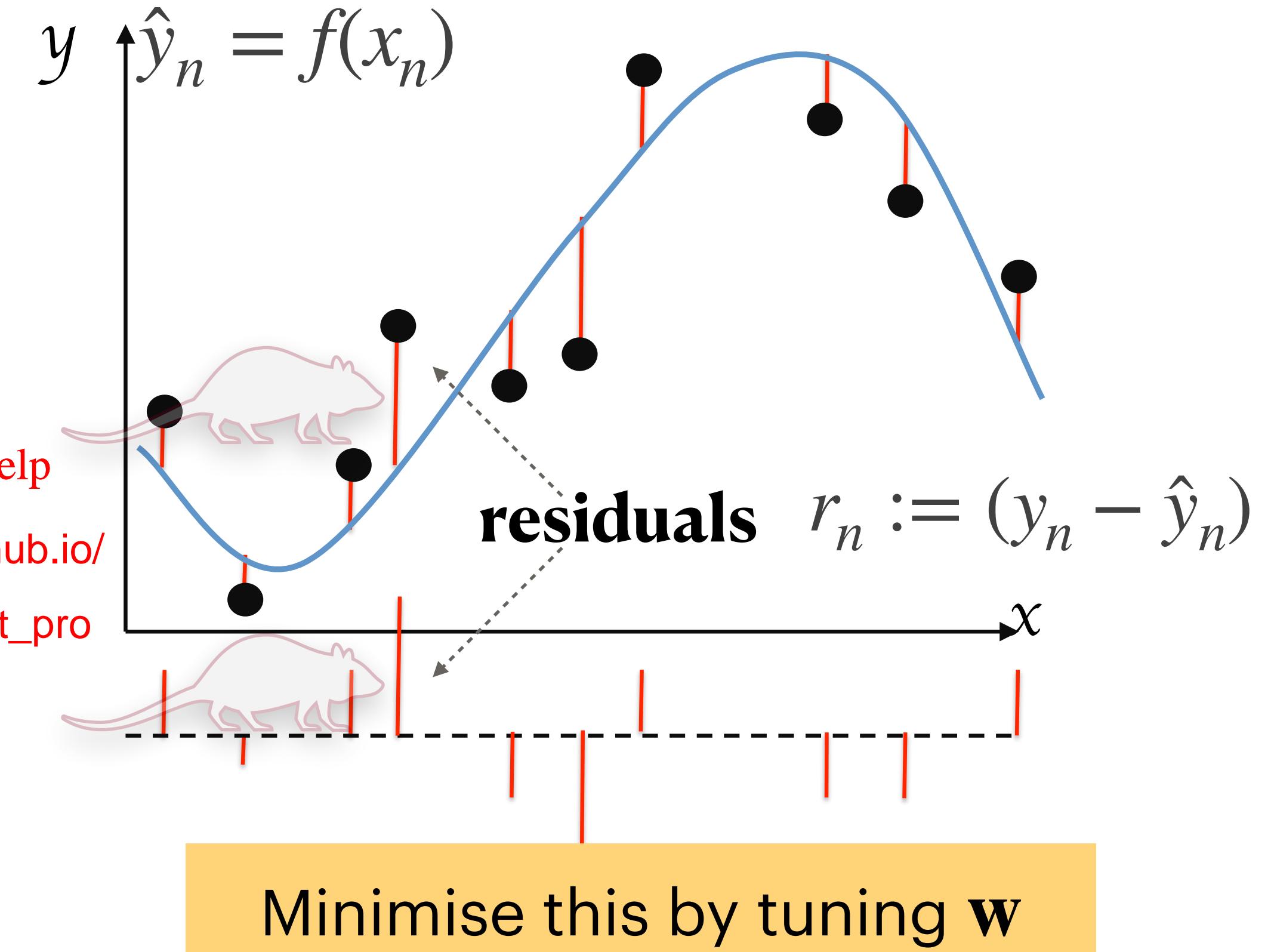
$$y_n = \hat{y}_n + \overbrace{(y_n - \hat{y}_n)}^{r_n} = f(x_n; \mathbf{w}) + r_n(\mathbf{w})$$

- Contribution to loss function:

$$l_n(\mathbf{w}) = (y_n - f(x_n; \mathbf{w}))^2 = r_n^2(\mathbf{w}).$$

- Average loss: $L(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N l_n(\mathbf{w})$

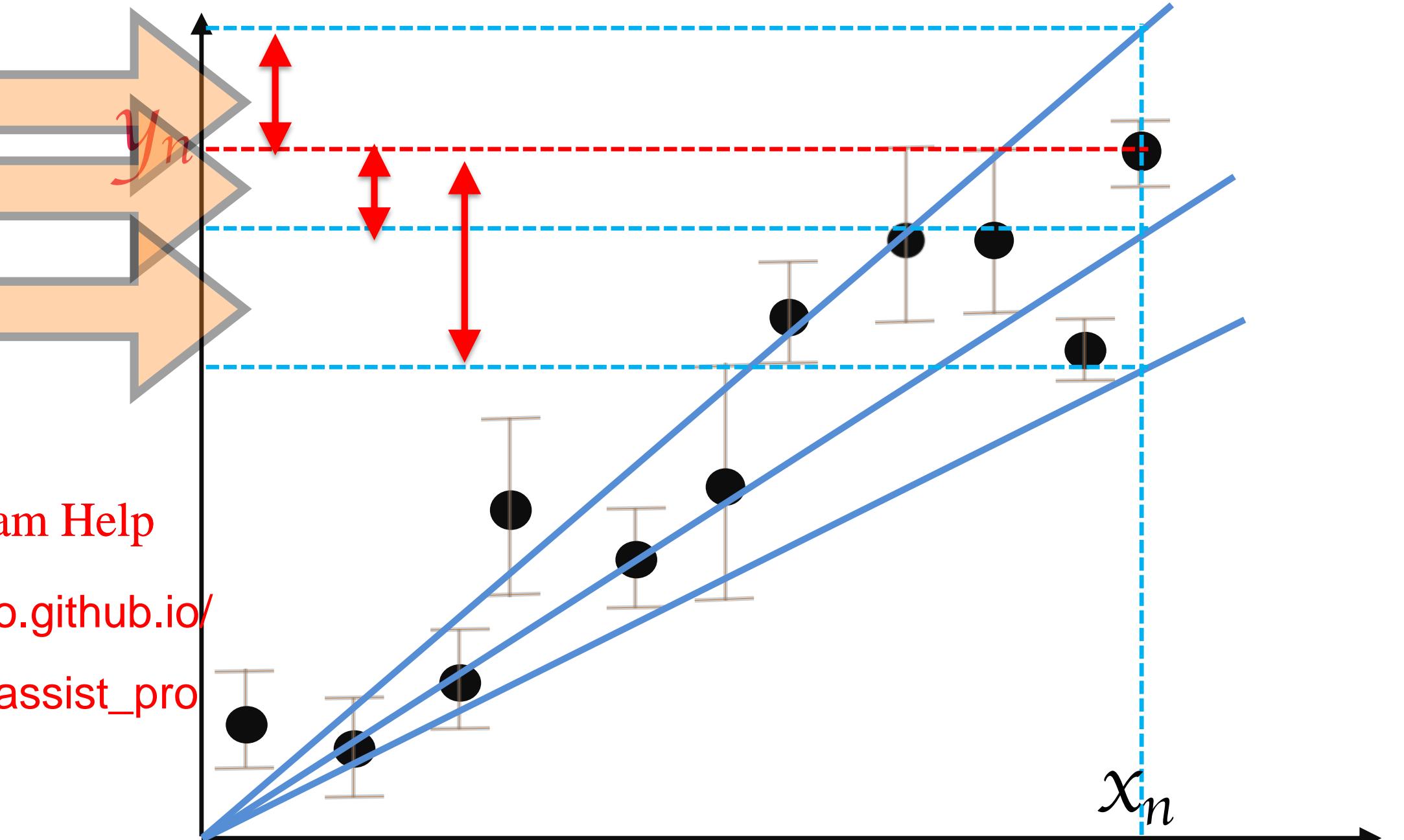
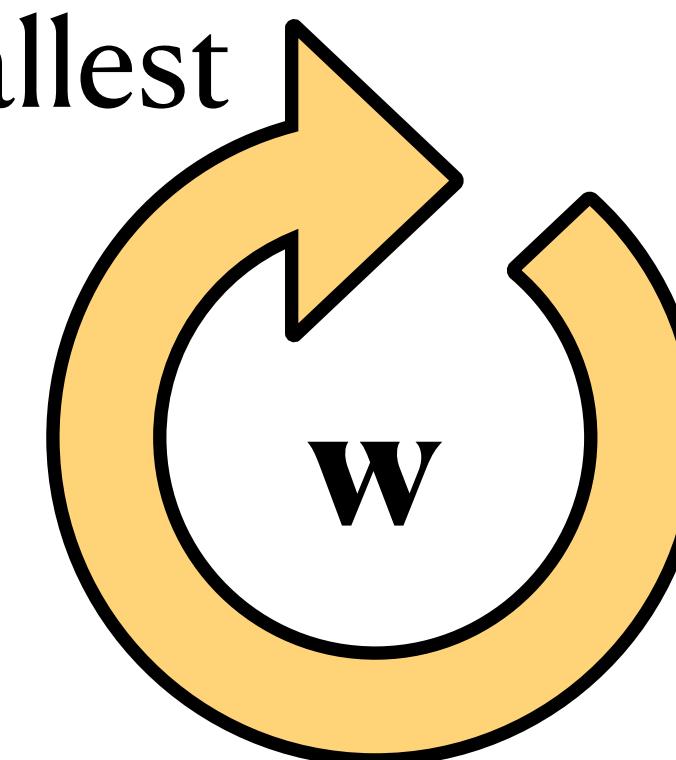
Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro



Choose weight for minimum loss

Example: find slope of straight line

- Fit $y = w x$ to data, slope w of the line is the weight/parameter to be learnt
- Loss L = sum of squares of residuals (red), for 3 possible choices of slope $\{w_1, w_2, w_3\}$: the 3 residuals for input x_n is shown
- Choose the slope that gives the smallest value from $\{L(w_1), L(w_2), L(w_3)\}$



$$l_n(w) = r_n^2 = (wx_n - y_n)^2$$

$$L(w) = \frac{1}{N} \sum_{n=1}^N l_n(w)$$

Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro

Loss function needs a distance: introducing the norm

Treat all data points as collective unit

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

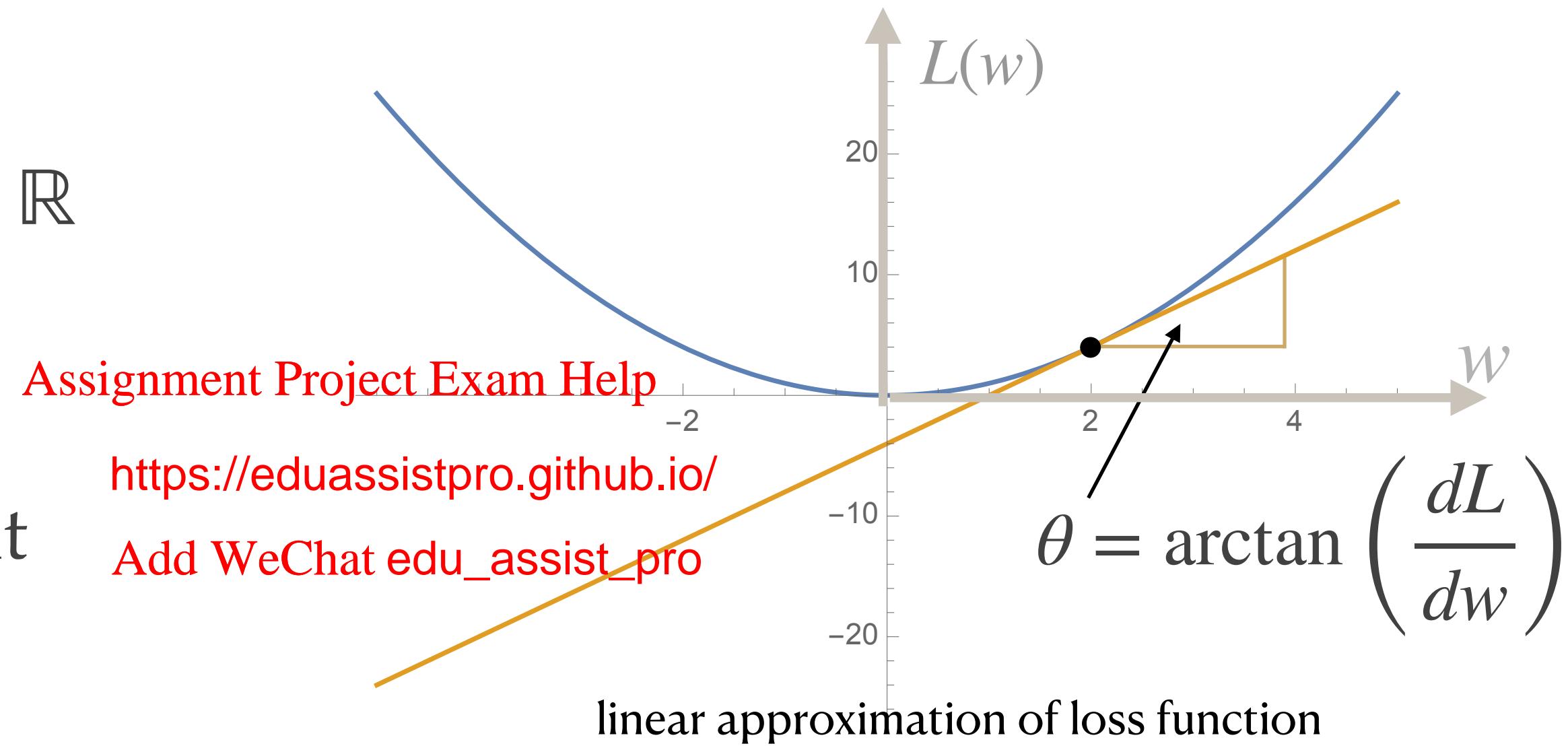
$$\mathbf{r} = \begin{pmatrix} r_1 \\ r_2 \\ \vdots \\ r_N \end{pmatrix}, \quad \|\mathbf{r}\|^2 = \sum_{n=1}^N r_n^2$$

Loss = (1/N) (length)² of N-dimensional residual vector

Update weights to reduce loss: gradient descent

Differentiable loss function: $l_n(w) = r_n^2 = (wx_n - y_n)^2$

- Slope can take any real value $w \in \mathbb{R}$
- Iterate $w^{(t)} \mapsto w^{(t+1)}$, $t = 0, 1, \dots$,
- Update weights $w^{(t)}$ to $w^{(t+1)}$ so that $L(w^{(t+1)}) < L(w^{(t)})$
- Change weights in the direction **opposite** to the slope of the loss function (we want to **reduce** the loss, hence **descent**)



$$w^{(t+1)} = w^{(t)} + \eta \left(\frac{dL(w^{(t)})}{dw} \right), \eta < 0$$

Linear Regression: solving for zero gradient of loss

$$(w^2x_1^2 - (w^2x_N^2 - 2wx_Ny_N + y_N^2))$$

Closed form solution exists: linear algebra

- $L(w) = (1/N)[(wx_1 - y_1)^2 + (wx_2 - y_2)^2 + \dots + (wx_N - y_N)^2]$

- Loss function is **quadratic in w**

- $L(w) = aw^2 + bw + c$

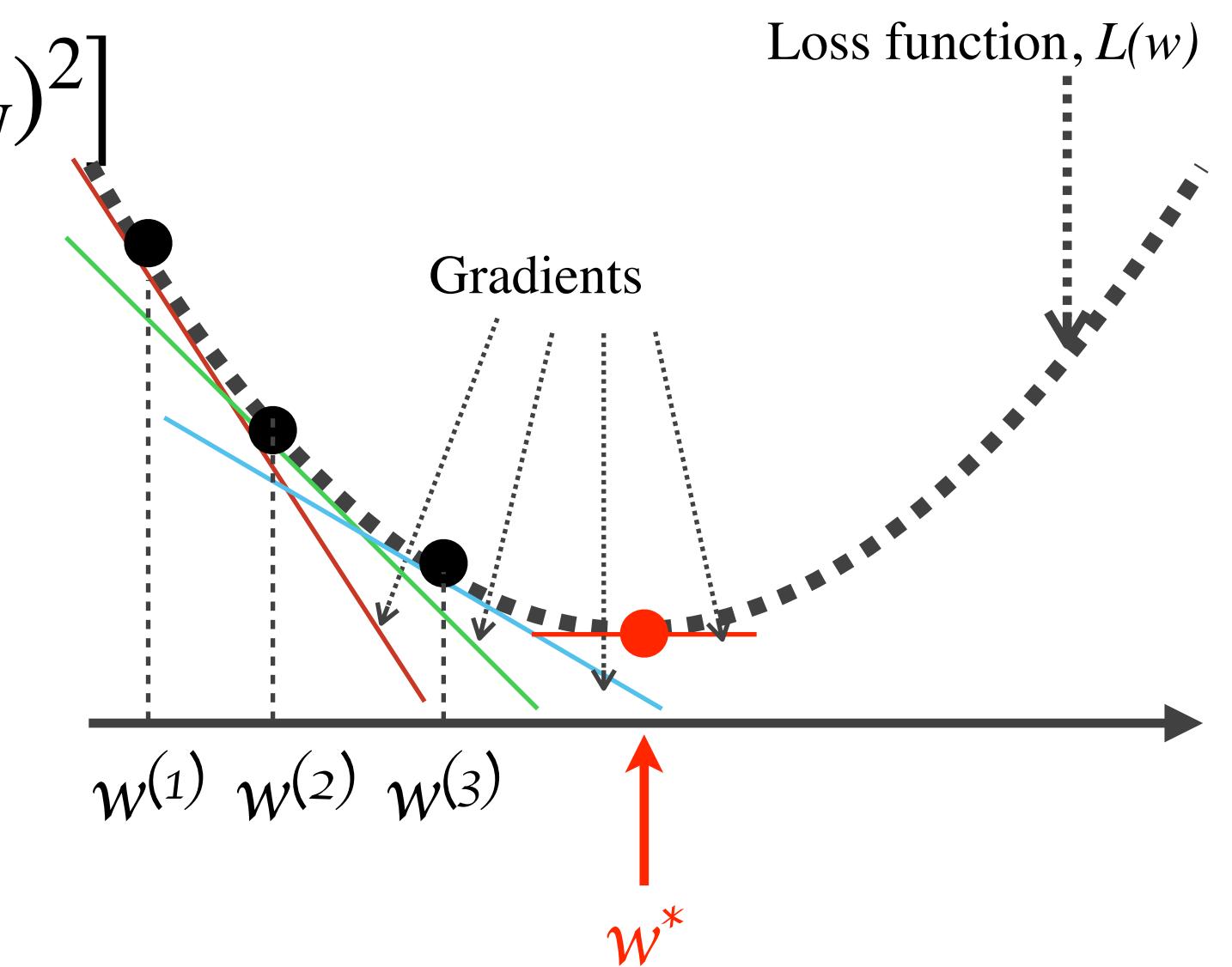
- Follow gradients until minimum reached

- Solution for weights: set **gradient = 0**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

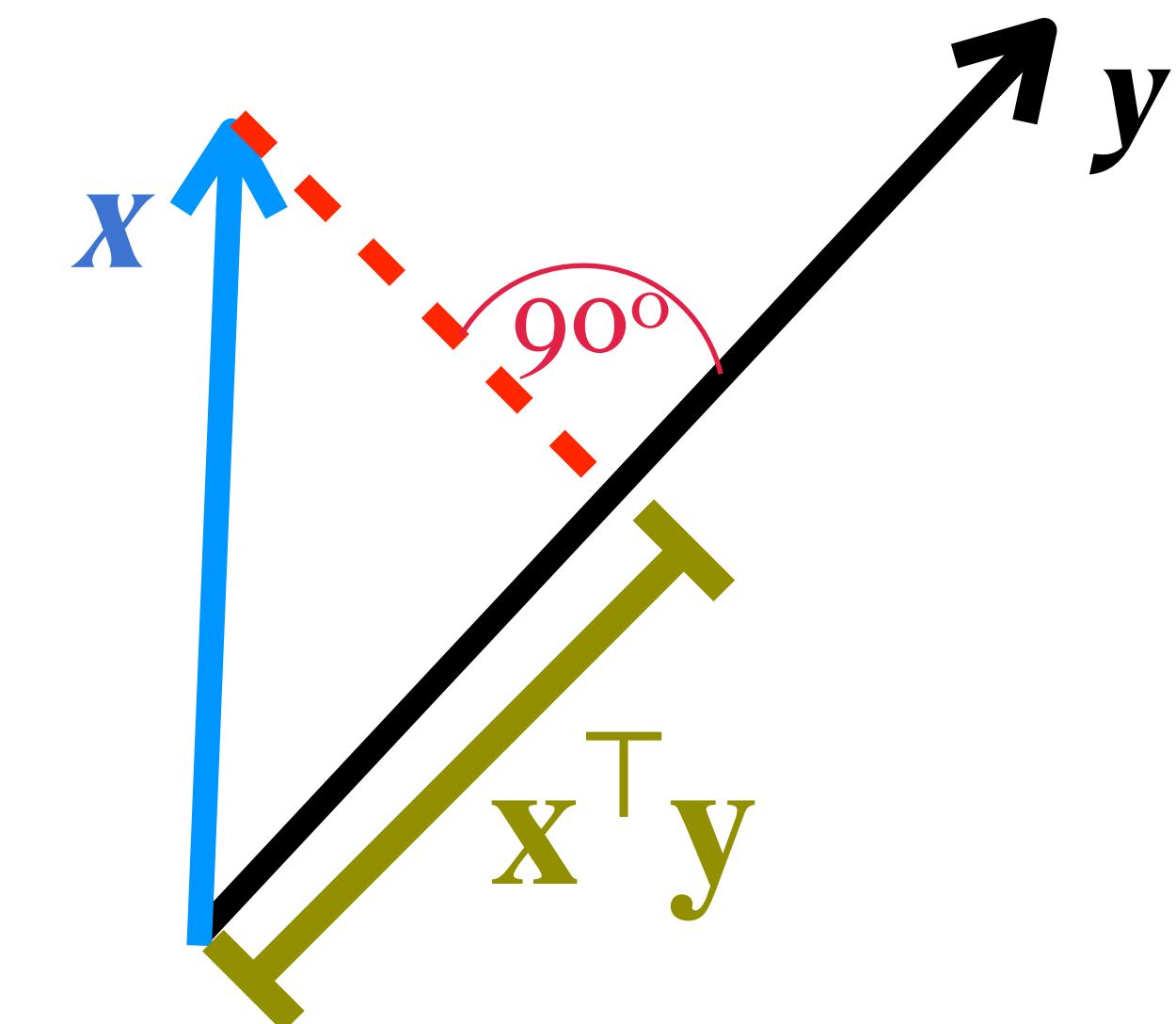


$$0 = \frac{\partial L(w)}{\partial w} \Big|_{w=w^*} \implies w^* = -b/(2a)$$

Exercise: differential calculus

Closed form solution to linear regression weights in terms of vector products

- $L(w) = (1/N) \left[(wx_1 - y_1)^2 + (wx_2 - y_2)^2 + \cdots + (wx_N - y_N)^2 \right]$
- **Exercise:** In $L(w) = aw^2 + bw + c$, show
Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro
- $a = (1/N)[x_1^2 + x_2^2 + \cdots + x_N^2]$
- $b = (-2/N)[x_1y_1 + x_2y_2 + \cdots + x_Ny_N]$
- $0 = \frac{\partial L(w)}{\partial w} \Bigg|_{w=w^*} \implies w^* = -b/(2a) = \frac{\mathbf{x}^\top \mathbf{y}}{\mathbf{x}^\top \mathbf{x}}$

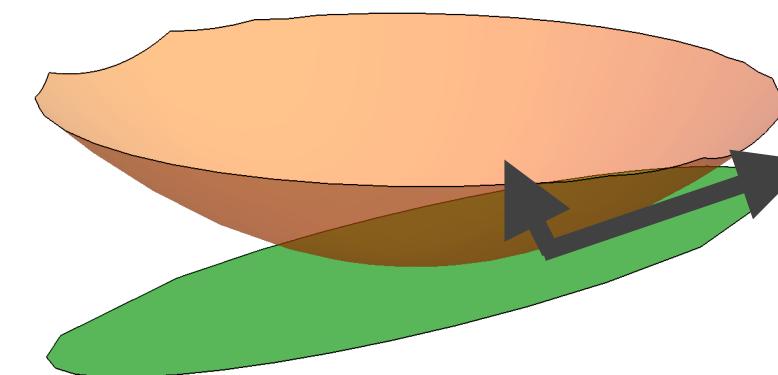


Reduce loss by gradient descent: optimisation

Same idea in higher dimensions (more adjustable weights)

$$\text{loss}_n = (y_n - (w_1 x_{n,1} + w_2 x_{n,2}))^2, \mathbf{x} = (x_1, x_2) \in \mathbb{R}^2$$

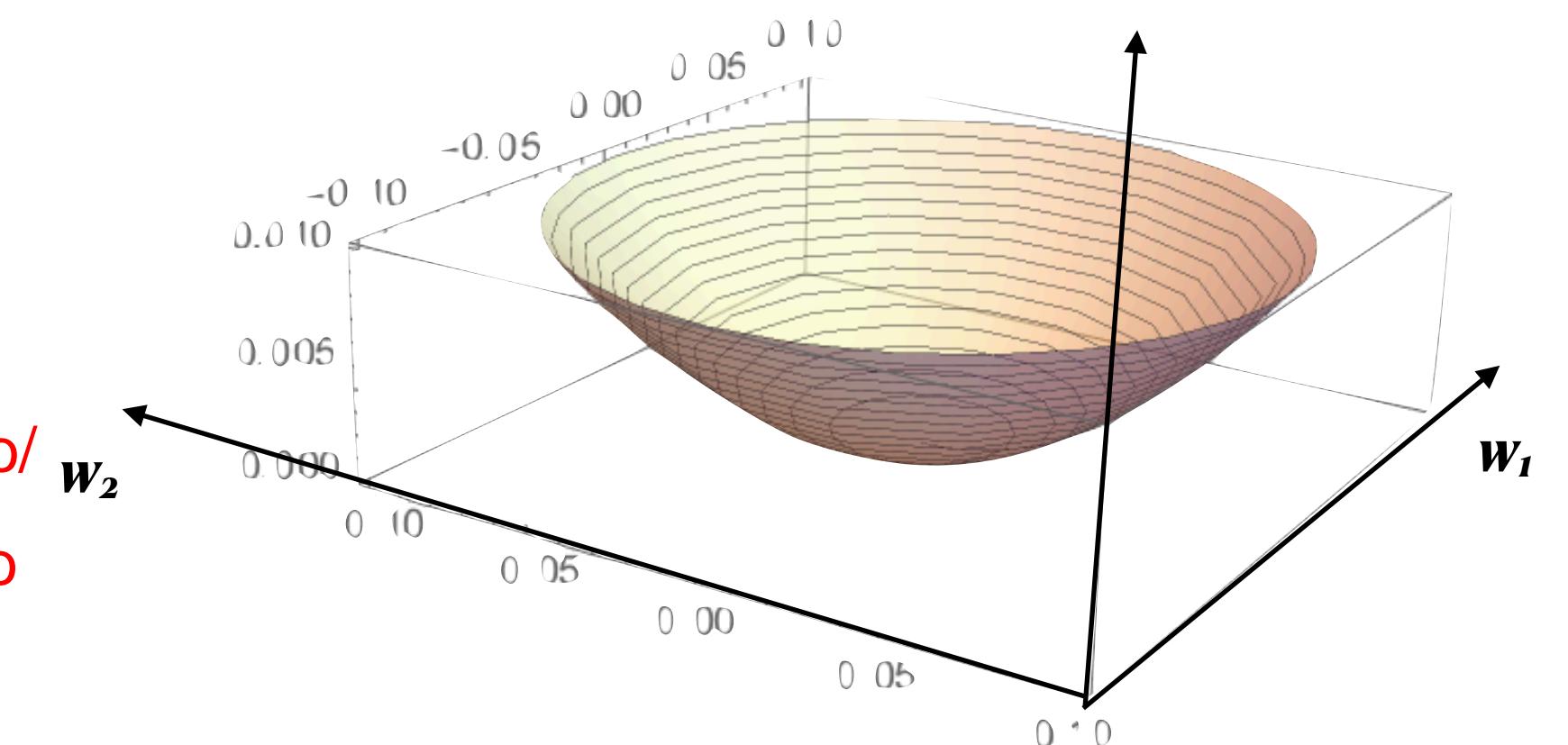
Evaluate partial derivatives (gradi
direction of weight updates



$$(\nabla_{\mathbf{w}} L)_i = \frac{\partial L}{\partial w_i}$$

Assignment Project Exam Help

se
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro



$$L(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}_n, y_n) \sim \mathcal{D}} \text{loss}_n(\mathbf{w}), \text{mean loss}$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} L$$

Automatic differentiation

No need to differentiate by hand, except to understand (lab exercises)

- Autograd, JAX – AD libraries in python

- $f: x \mapsto f(x); \frac{dy}{dx} = \lim_{\delta x \rightarrow 0} \frac{f(x + \delta x) - f(x)}{\delta x} = \lim_{\delta x \rightarrow 0} \frac{f(x + \delta x) - f(x - \delta x)}{2 \cdot \delta x}$
Assignment Project Exam Help

- Product: $y(x) = f(x) \cdot g(x); y'(x) = f'(x) \cdot g$ <https://eduassistpro.github.io/>

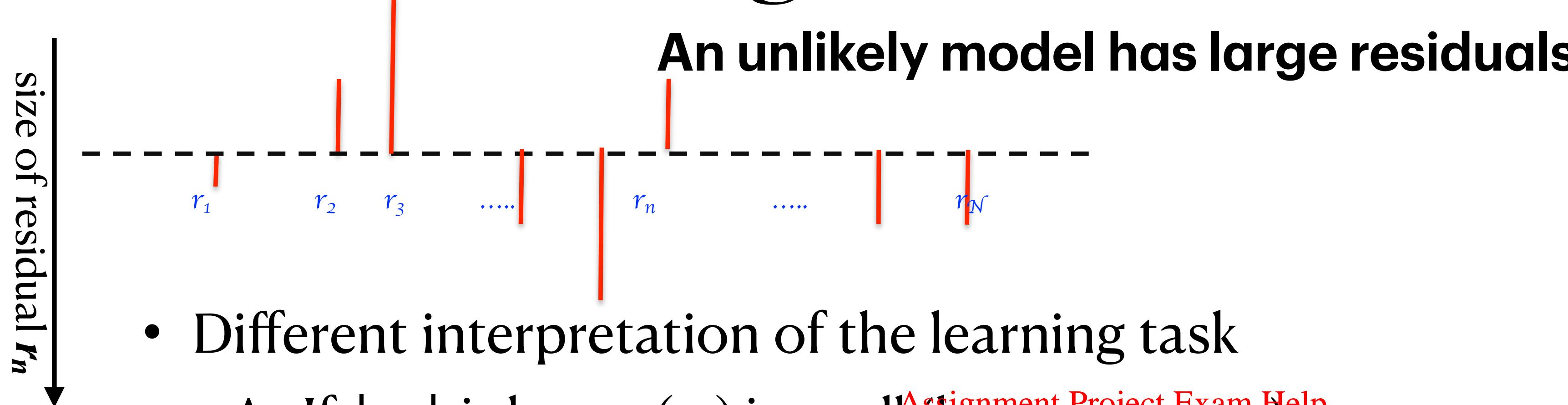
- Quotient: $y(x) = f(x)/g(x); y'(x) = [f'(x) \cdot g(x) - f(x) \cdot g'(x)]/g(x)^2$
Add WeChat edu_assist_pro

- Composition: $y(x) = f(g(x)), \frac{dy}{dx} = \frac{df}{dg} \frac{dg}{dx}$

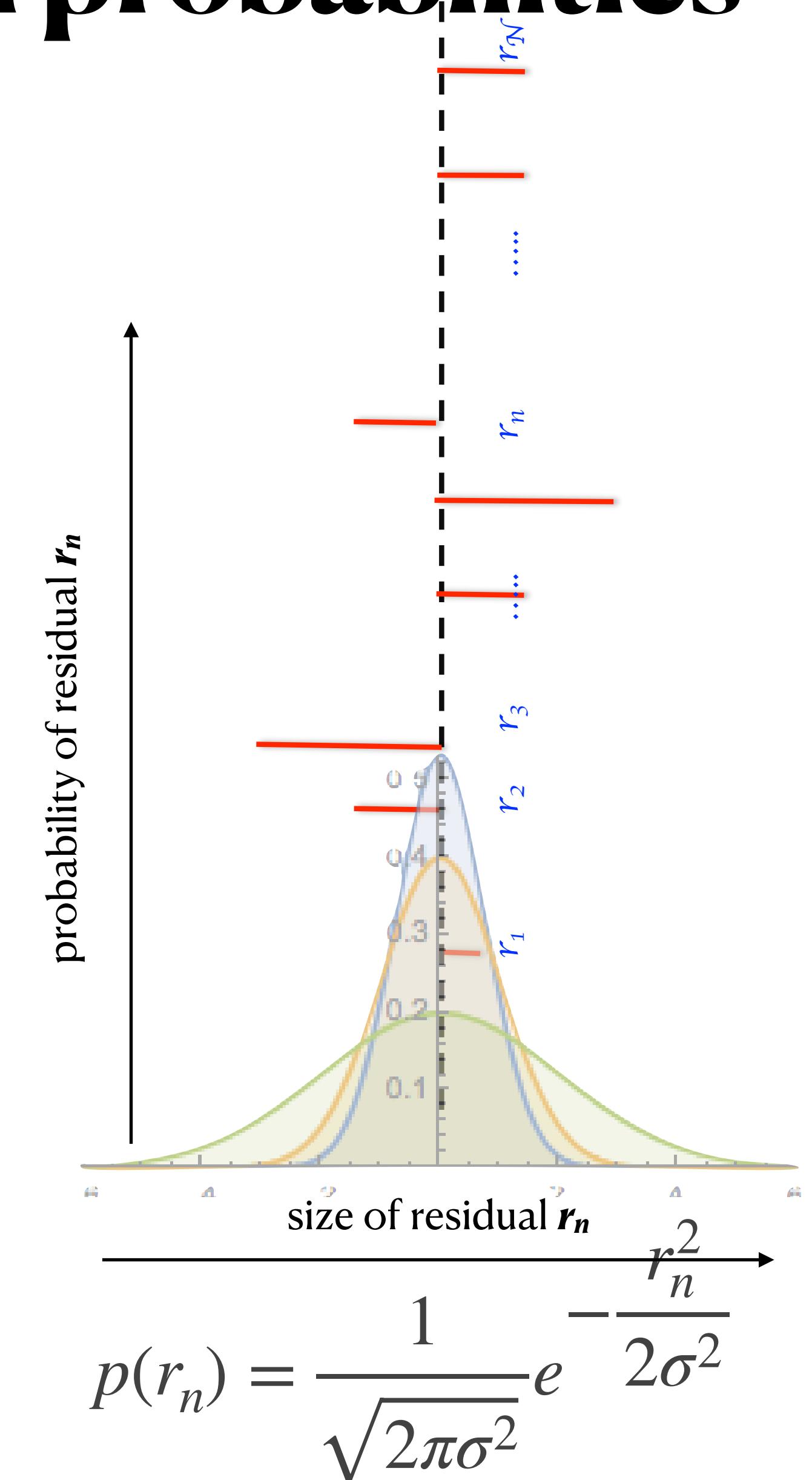
- Program $f: x \mapsto f(x)$ forward mode AD (f, Df): $x \mapsto (f(x), f'(x))$ e.g., $\text{sq}: x \mapsto (x^2, 2x)$

- Taylor series: $f(x_0 + \delta x) = f(x_0) + \delta x \left. \frac{df}{dx} \right|_{x_0} + \frac{1}{2} (\delta x)^2 \left. \frac{d^2 f}{dx^2} \right|_{x_0} + \dots$, note $(\delta x)^2 \ll \delta x$

Later: view large distances as small probabilities

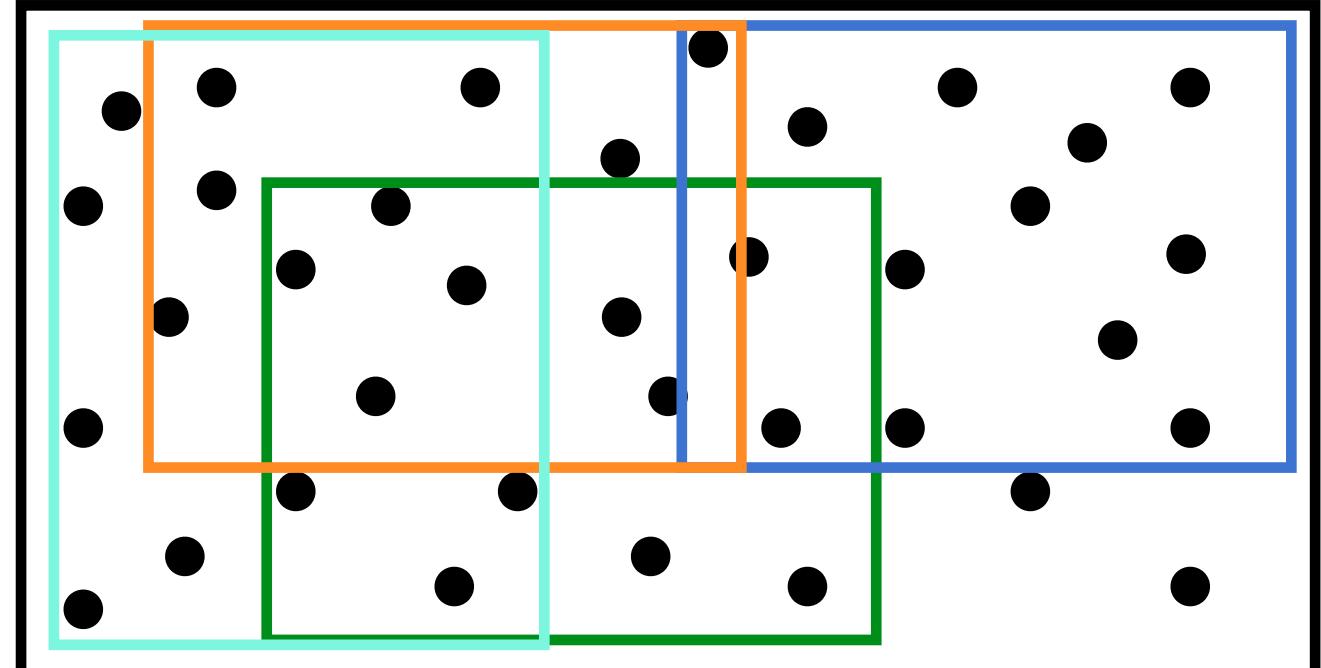


- Different interpretation of the learning task
 - A. If $|r_n|$ is large $p(r_n)$ is small (Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro)
 - B. Reducing $|r_n|$ makes $p(r_n)$ large (probability)
- Probability of what? $f(\cdot; \mathbf{w})$ evaluated on data
 $\mathcal{D} := \{(x_n, y_n)\}_{n=1, \dots, N}$, called model **likelihood**
- Maximum likelihood estimation: estimation of weights to achieve objective of large likelihood



Learning or memorising?

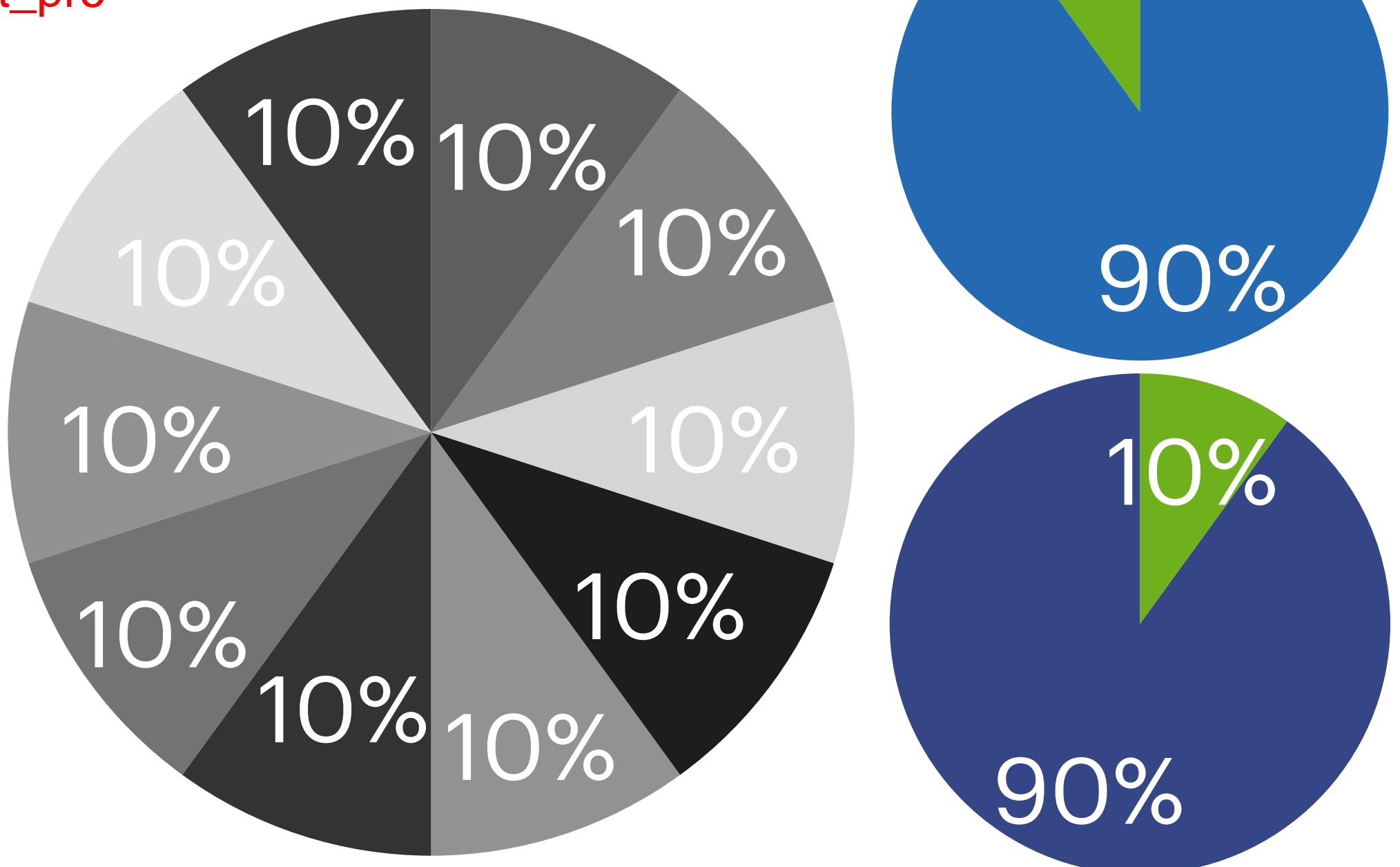
Evaluating ML models: Cross-validation



Coloured boxes: possible training sets

- Evaluate model performance on test data **(generalisation)**
- Different training sets can lead to model parameters with different predictions, hence different residuals
- Characterise model on **distribution** of residuals trained on different subsets of available training data
- **Cross validation**

Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat edu_assist_pro



Learning or memorising? Bias-variance tradeoff

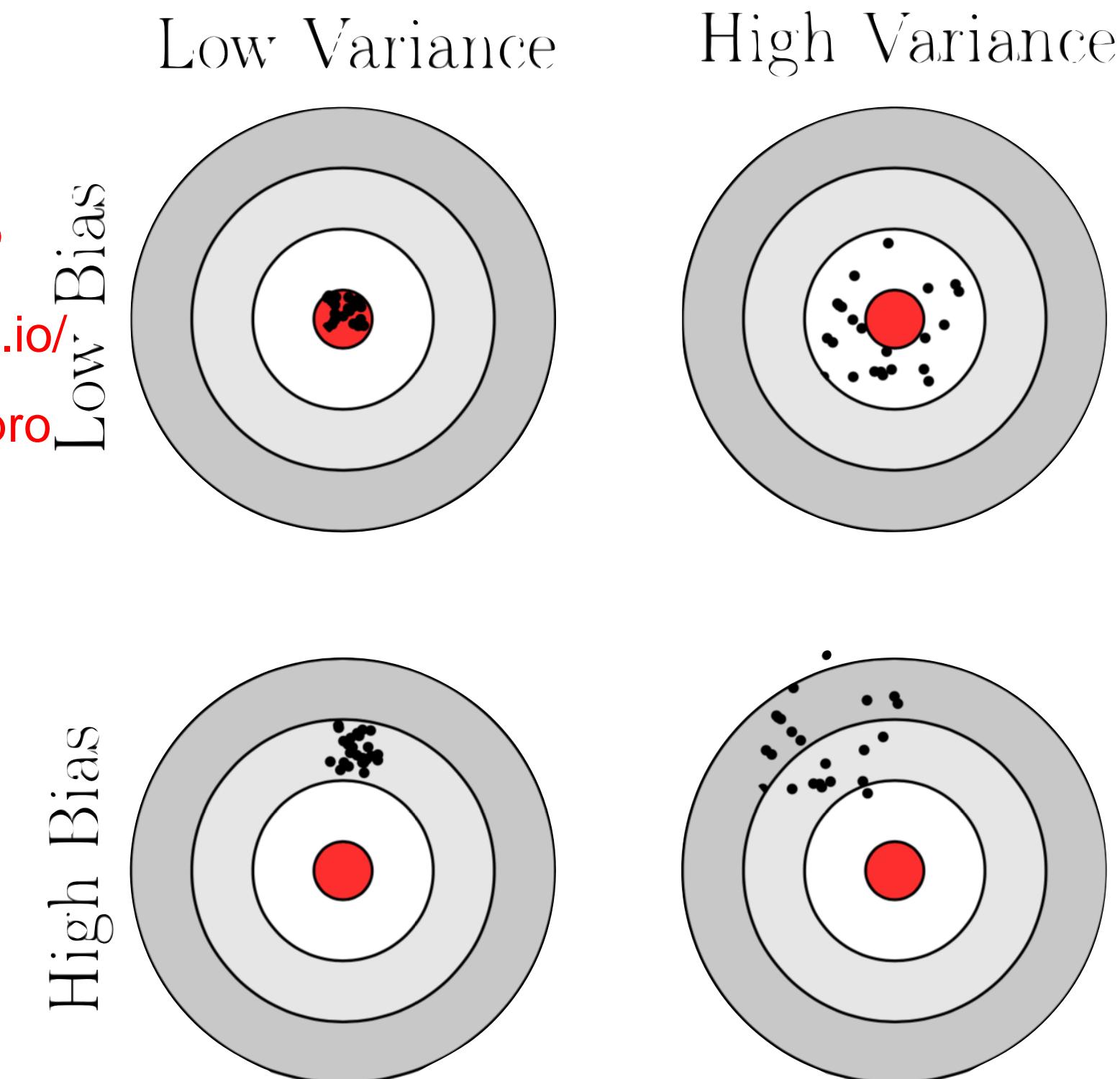
Criteria for evaluation of ML models

- Different training sets can lead to different models with different model predictions and different residuals
- Bias: Deviation of average of predictions from truth
- Variance: variability of model predictions on different parts of test set
- Trade-off

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro



Curve-fitting,
no interpretation

Assignment Project Exam Help

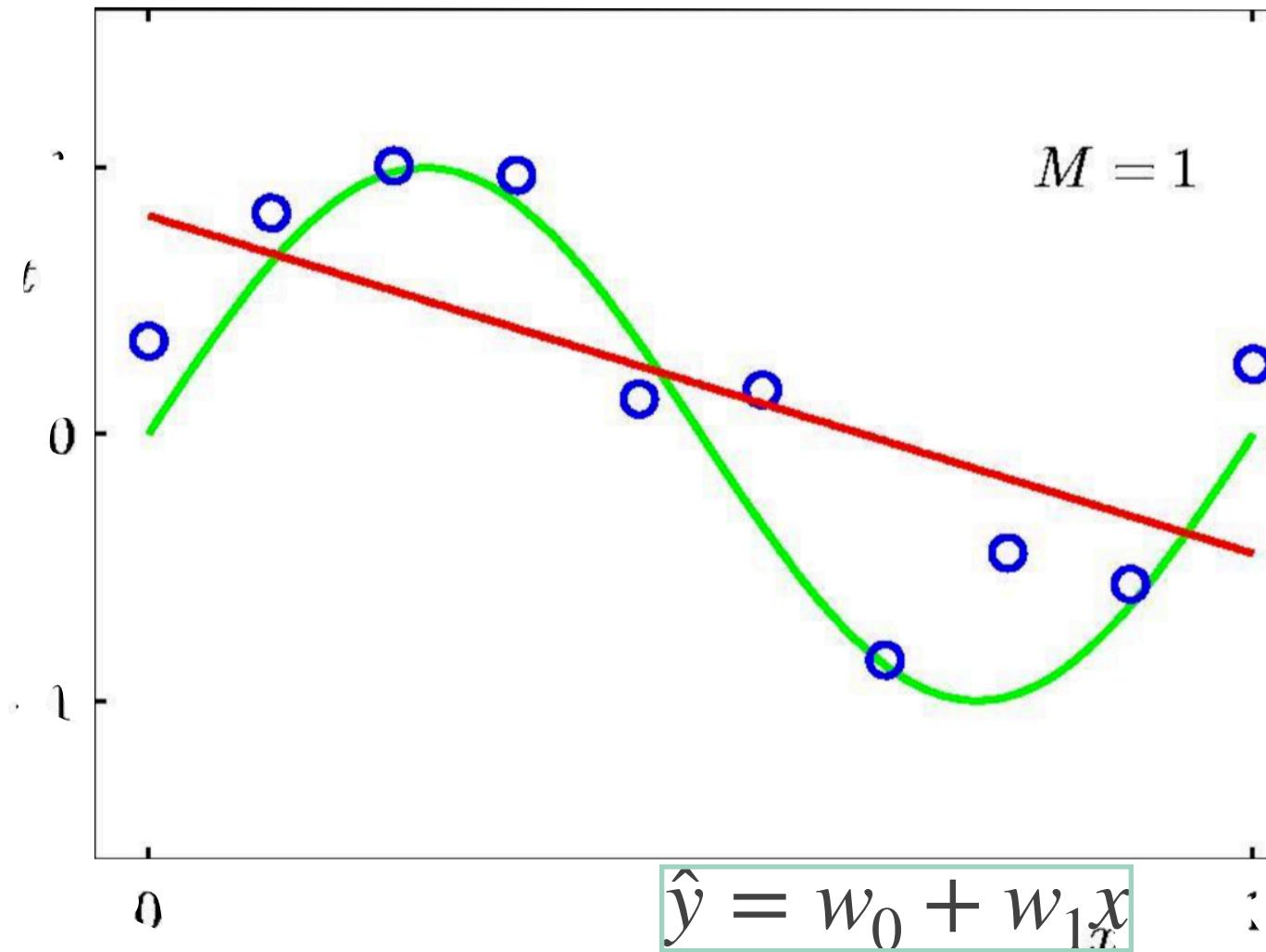
<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

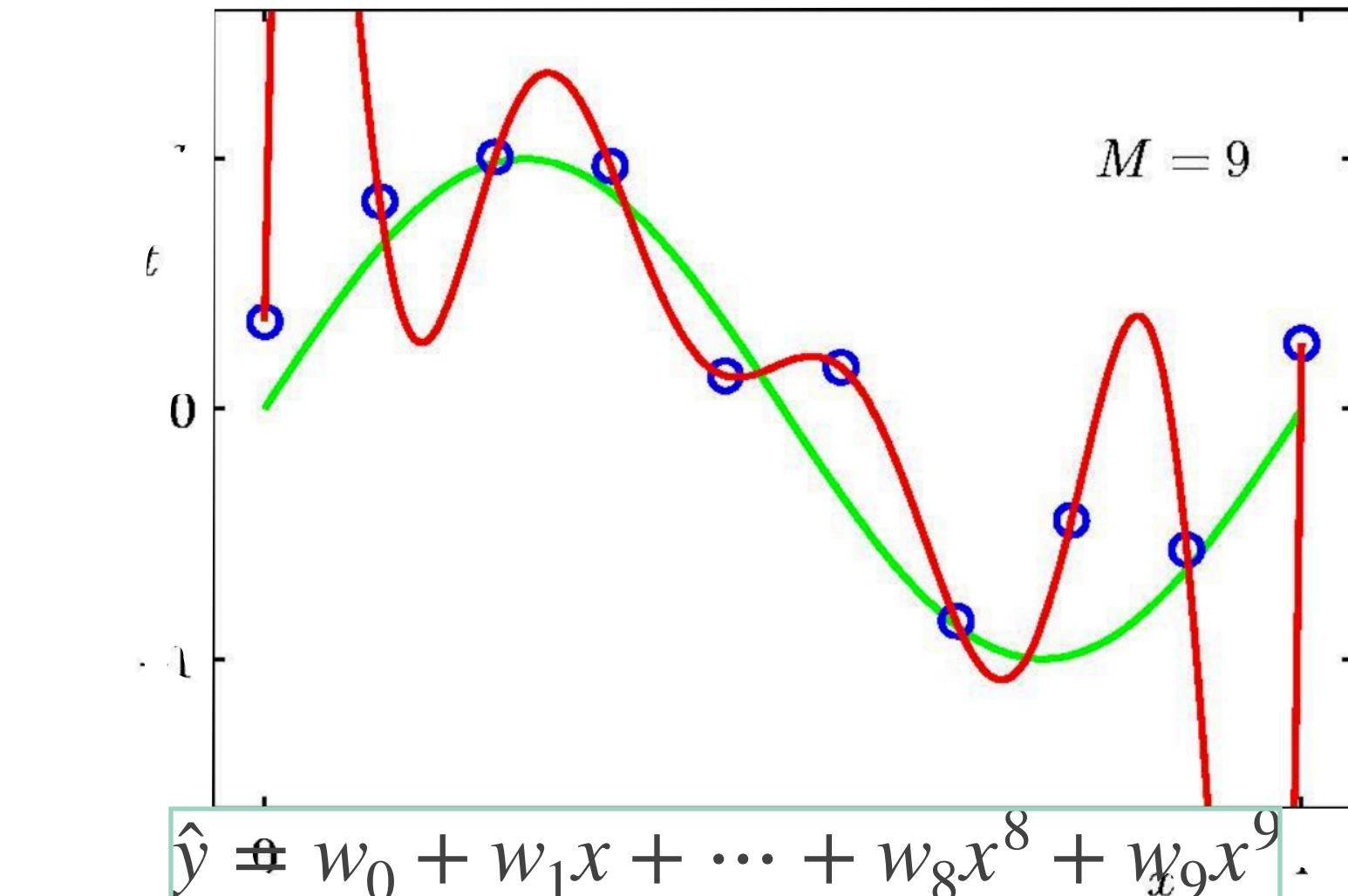


Polynomial fits with high degrees tend to overfit

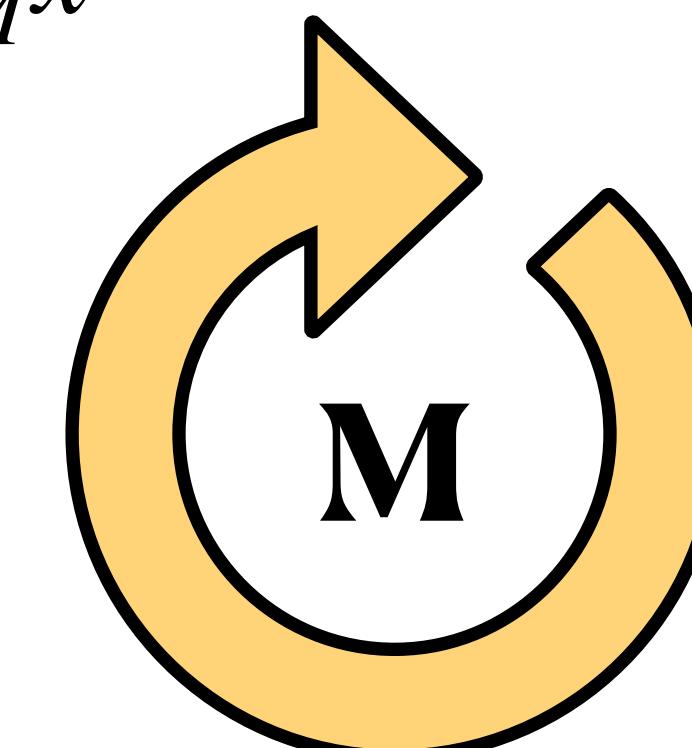
Overfitting: good on training set, poor on test set



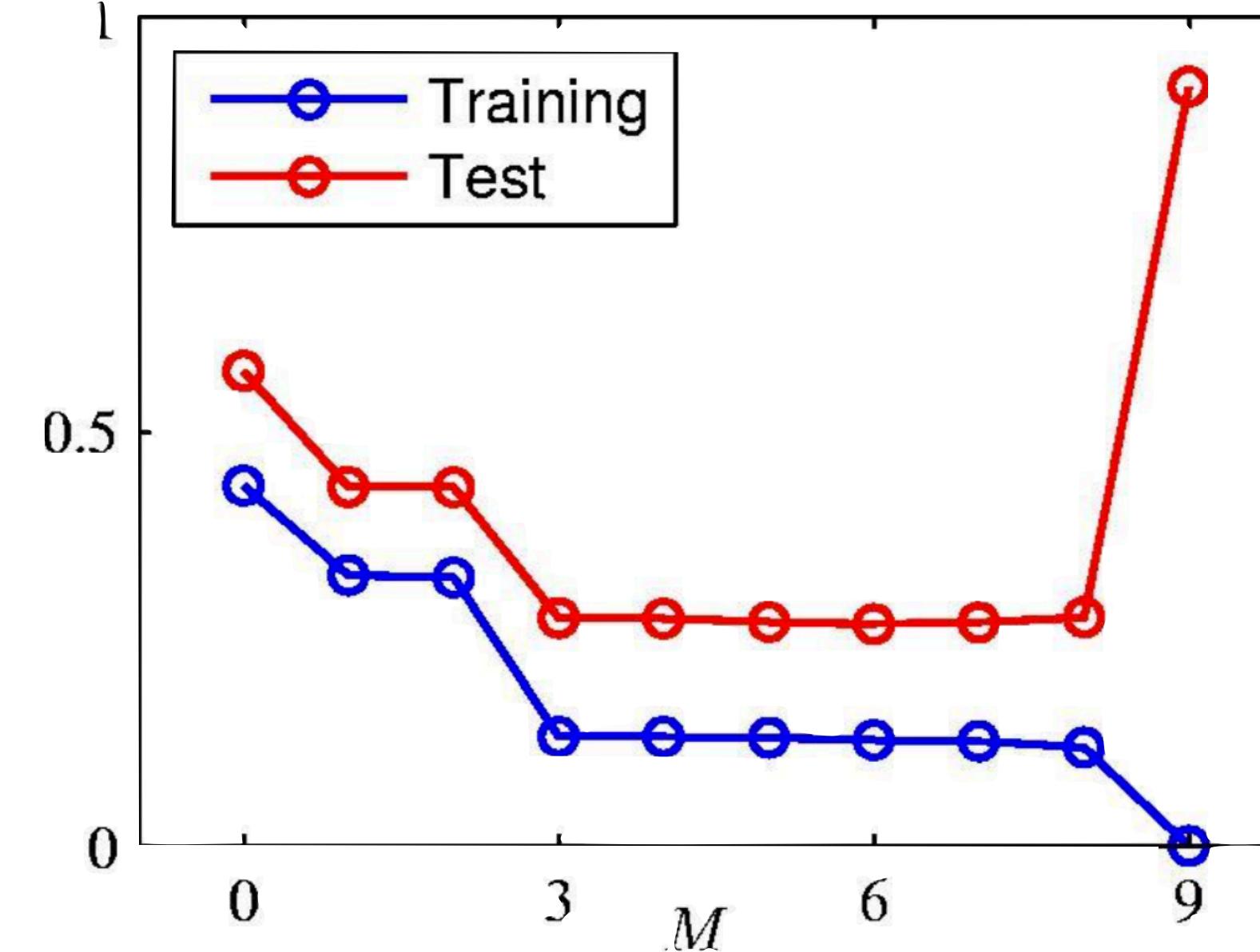
Assignment Project Exam Help
<https://eduassistpro.github.io/>
Add WeChat `edu_assist_pro`



$$\hat{y}(x; \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M$$



From C Bishop, PRML



Weights learned by minimising loss

Polynomial models of high degree have large weights

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu_assist_pro

From C Bishop, PRML

Regularisation: Penalty for model complexity

Complex models are believed to overfit

$$\text{LOSS}_{\mathcal{D}}(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}_n, y_n) \sim \mathcal{D}} (y_n - \hat{y}(\mathbf{x}_n; \mathbf{w}))^2$$

Assignment Project Exam Help

<https://eduassistpro.github.io/>

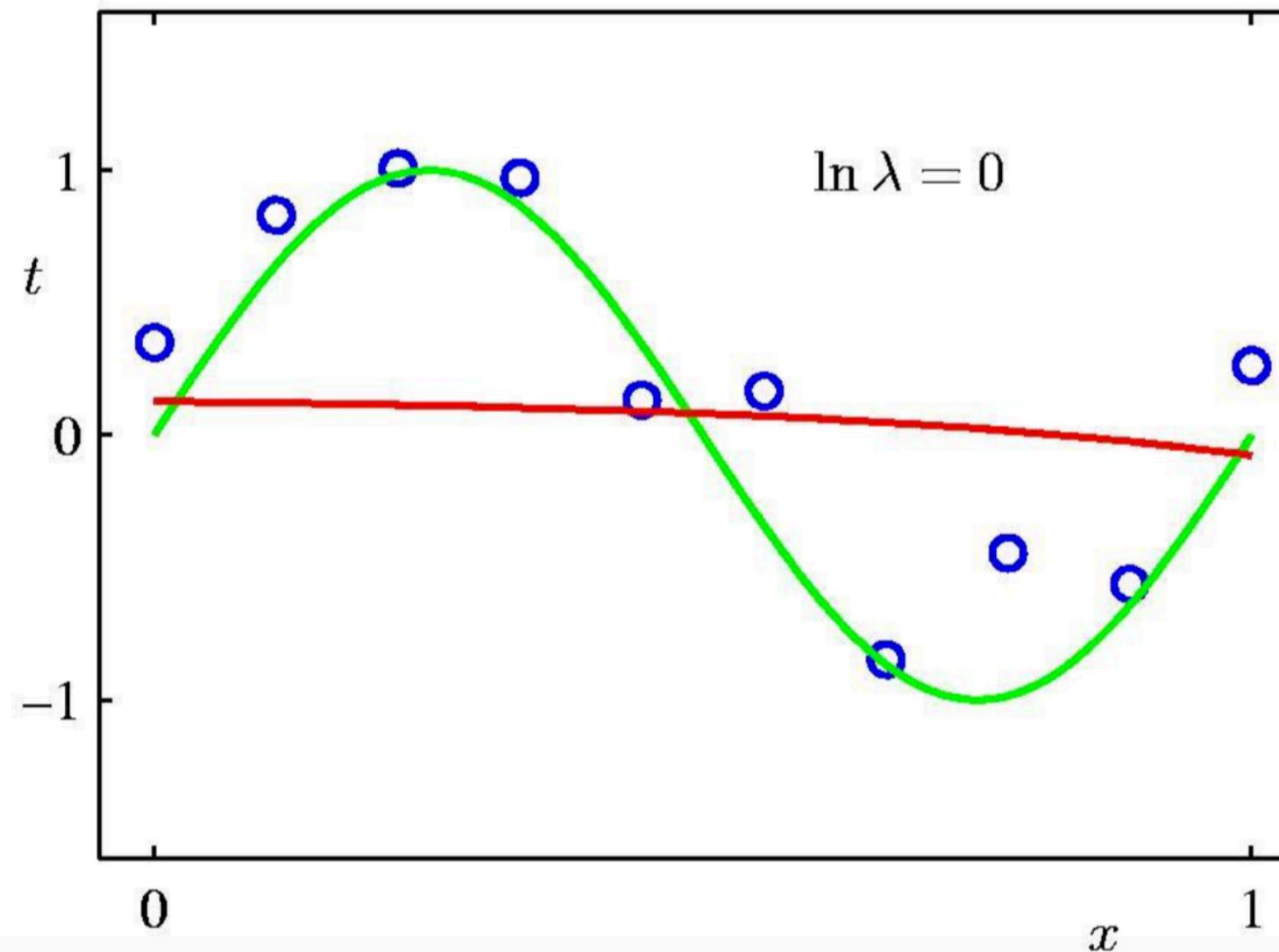
$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \text{LOSS}_{\mathcal{D}} + \lambda \|\mathbf{w}\|^2$$

Minimise a combination of two factors

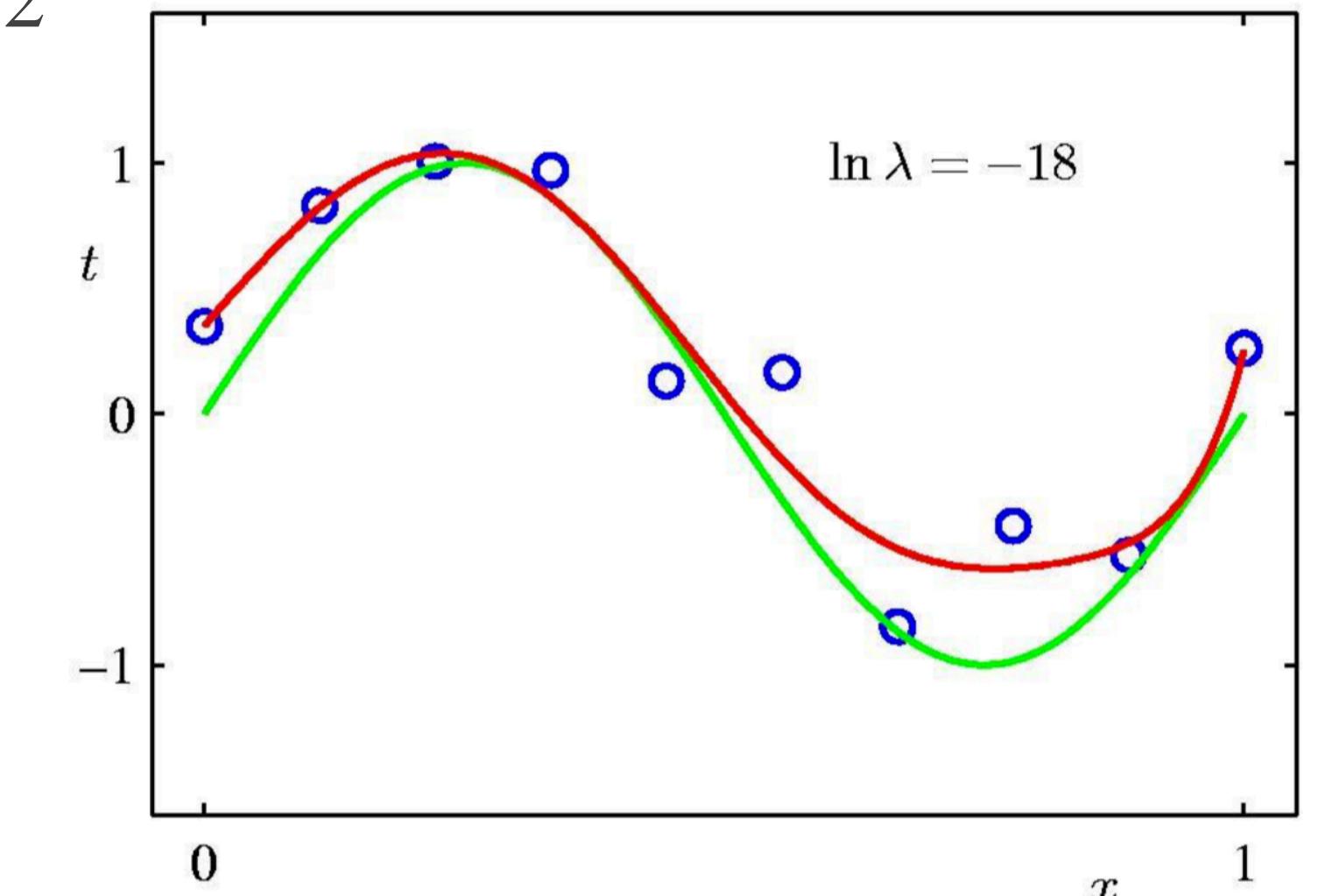
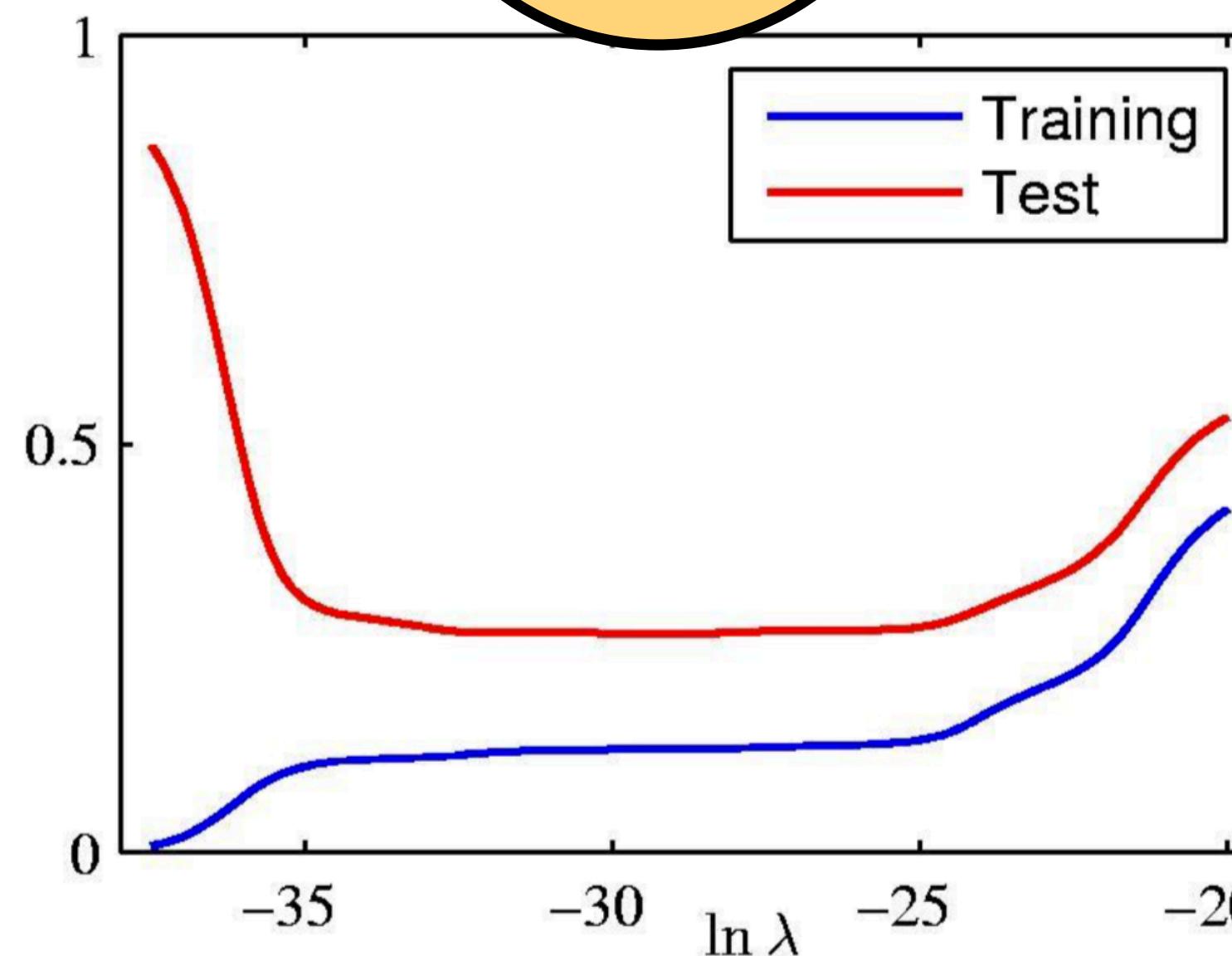
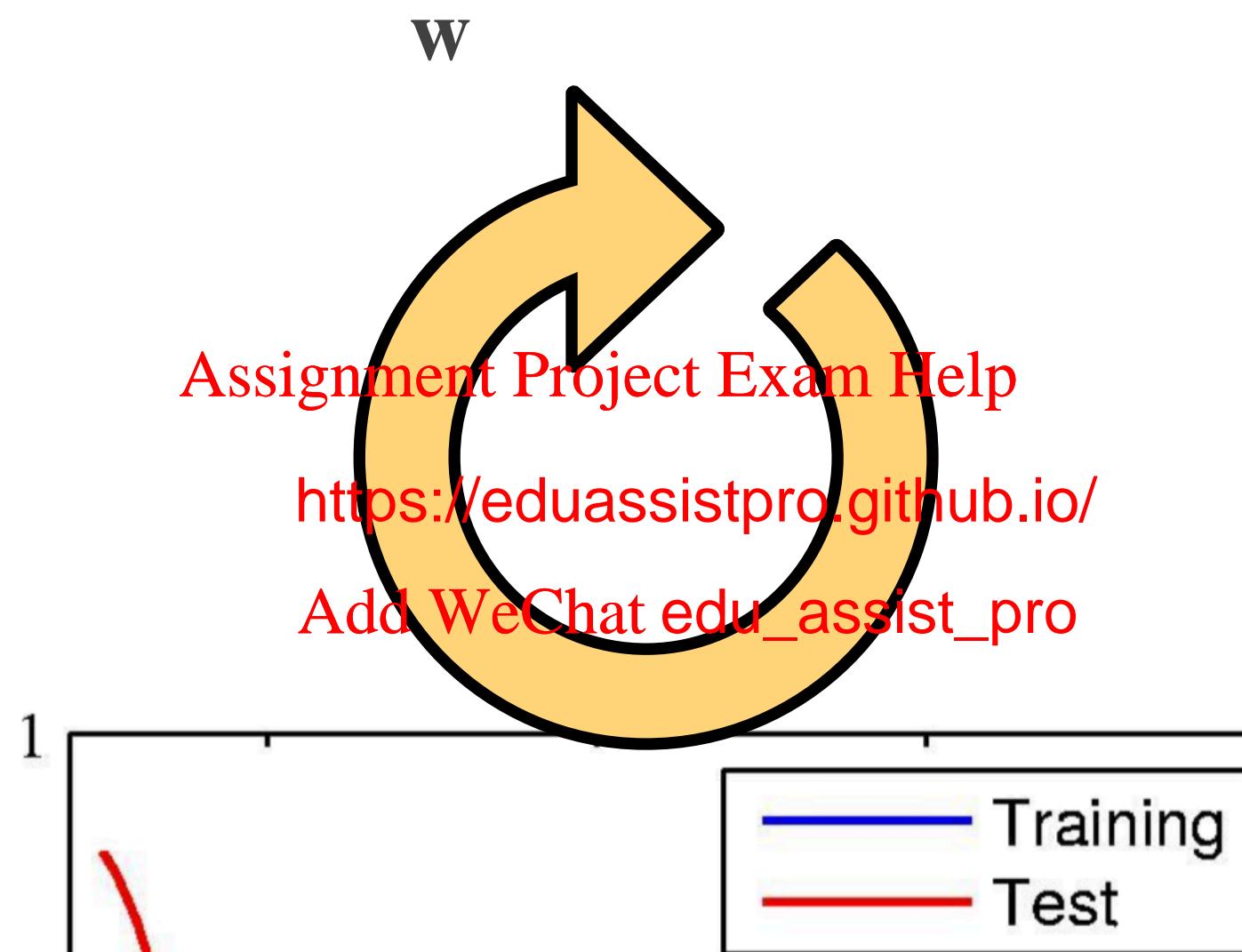
1. – mismatch of model prediction to labelled data (Loss)
2. – data-independent term that depends on model alone (**regularisation**)

Relative penalties for loss and weight length

How big should model penalty λ be?



$$\mathbf{w}^* = \arg \min \text{Loss}_{\mathcal{D}}(\mathbf{w}) + \lambda \|\mathbf{w}\|^2$$



*Read Section 1.1, p.4
C Bishop, PRML*

Supervised Learning: reduce loss

Summary

- Predict and correct using weight-adjustable functions
- Optimise weights using loss function – learning as optimisation
 - Assignment Project Exam Help
- Interpret weight spaces in terms of <https://eduassistpro.github.io/> the framework of linear algebra
- Interpret distances and loss metrics in terms of probability theory
 - Add WeChat [edu_assist_pro](#)
- Evaluate performance in terms of generalisation (bias/variance)
- Introduce regularisation for generalisation