

Parameters vs hyperparameters

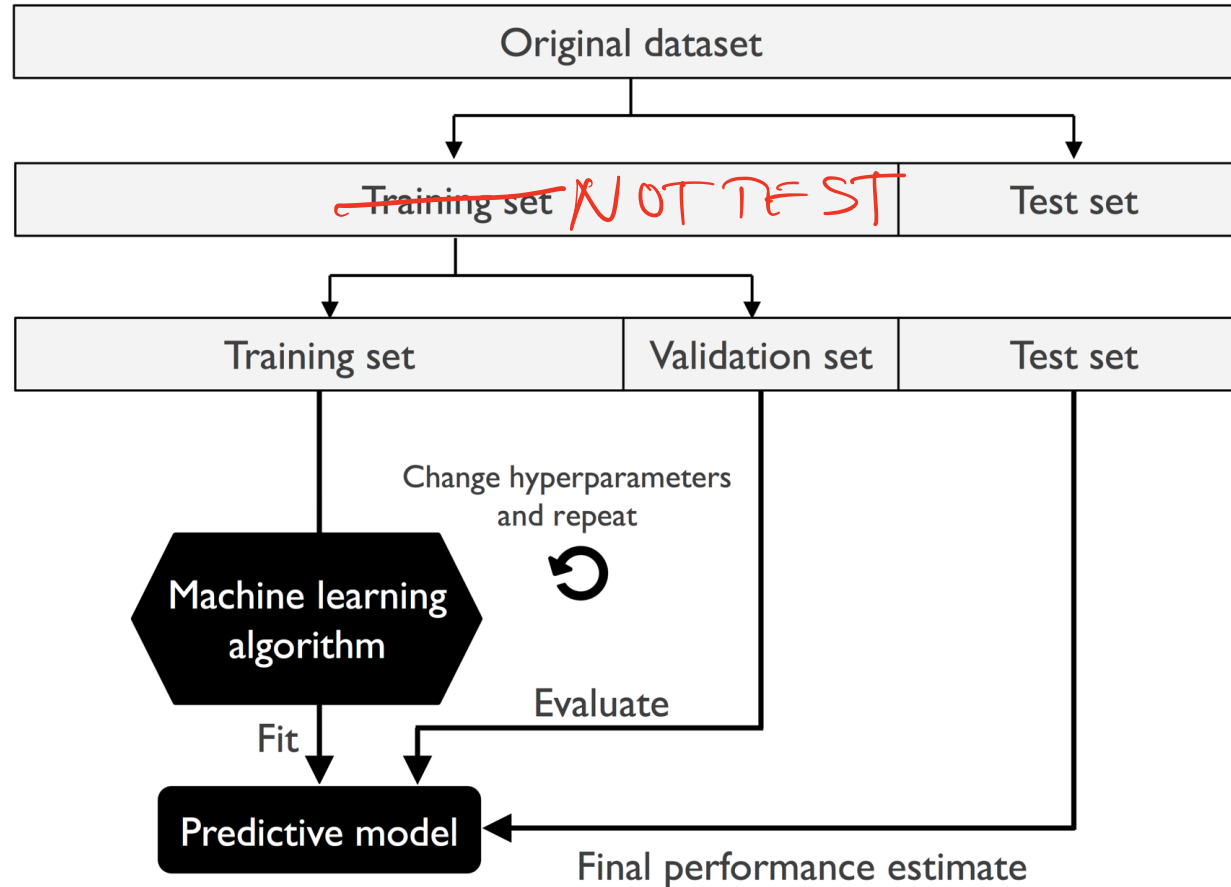
$$\hat{y} = f(\vec{x}) = 3x_1^2 + 2x_2$$

$\{$ parameters \equiv learned during
(weights) training to
(coefficients) minimize loss

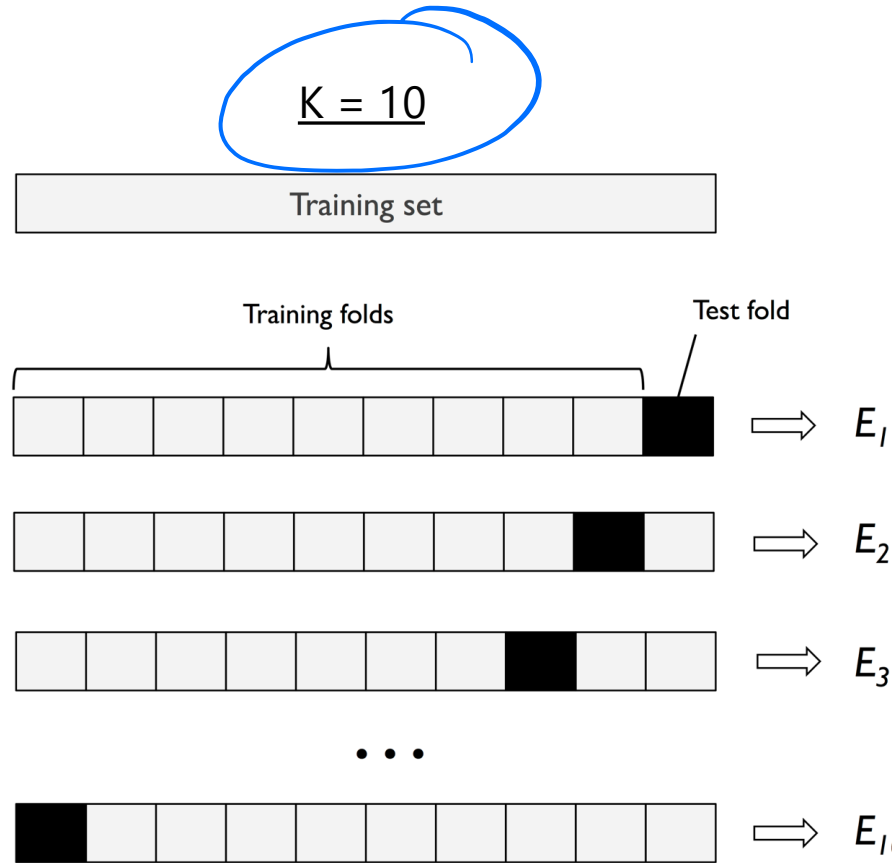
hyperparameters \equiv knobs that were turned to
arrive at $f(\vec{x})$

- ↳ # features
- ↳ loss function
- ↳ model type

Validation and testing workflow



K-fold cross validation



1) Choose hyperparameters

2) Run CV

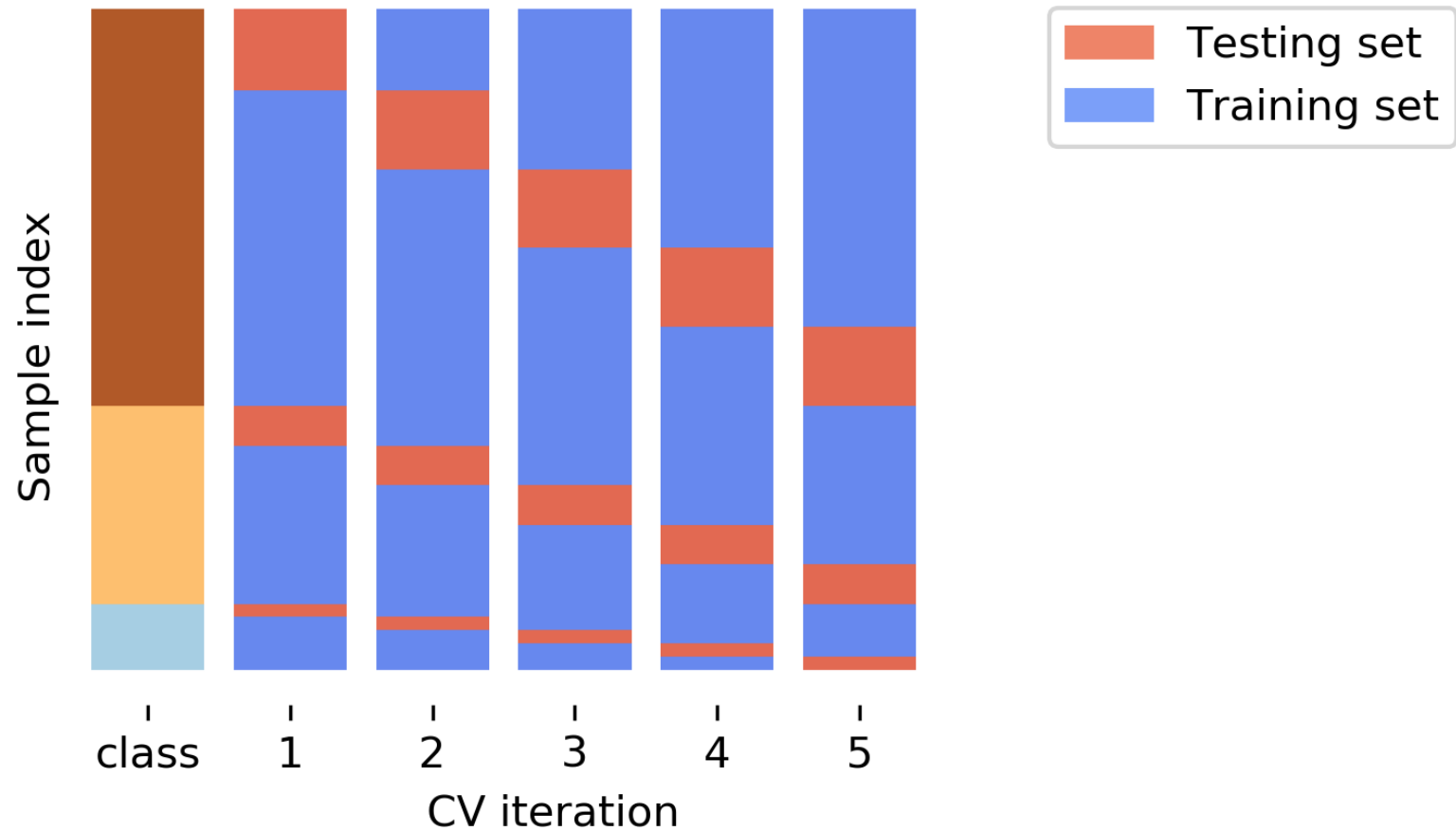
3) Average scores

$\uparrow K \Rightarrow$ more robust
estimate of
validation
performance

$\uparrow K \Rightarrow \uparrow$ cost

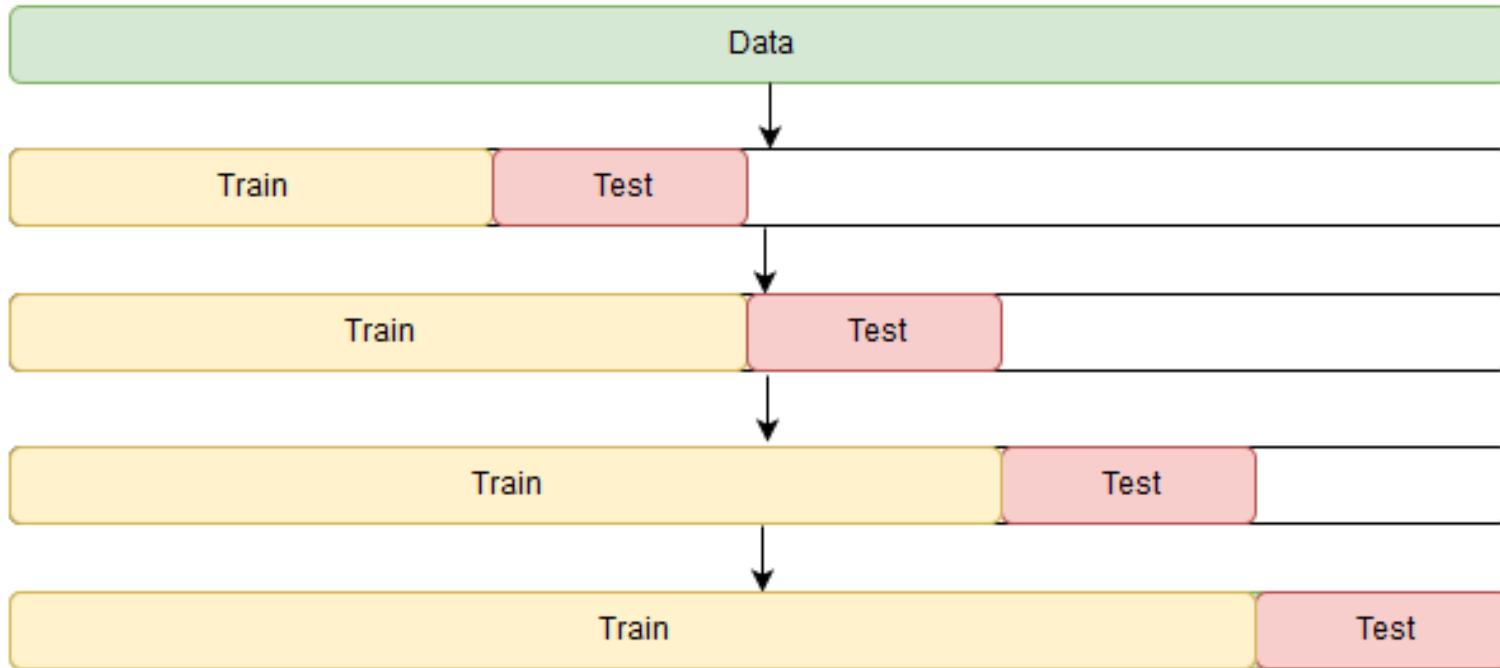
$$E = \frac{1}{10} \sum_{i=1}^{10} E_i$$

Class imbalance: stratified cross validation



Time series: "Rolling validation"

$t_0, t_1, t_2, \dots, t_f$



Understanding data validation: IID vs OOD

IID \Rightarrow independent, identical distribution
(sample randomly from same data
as training)

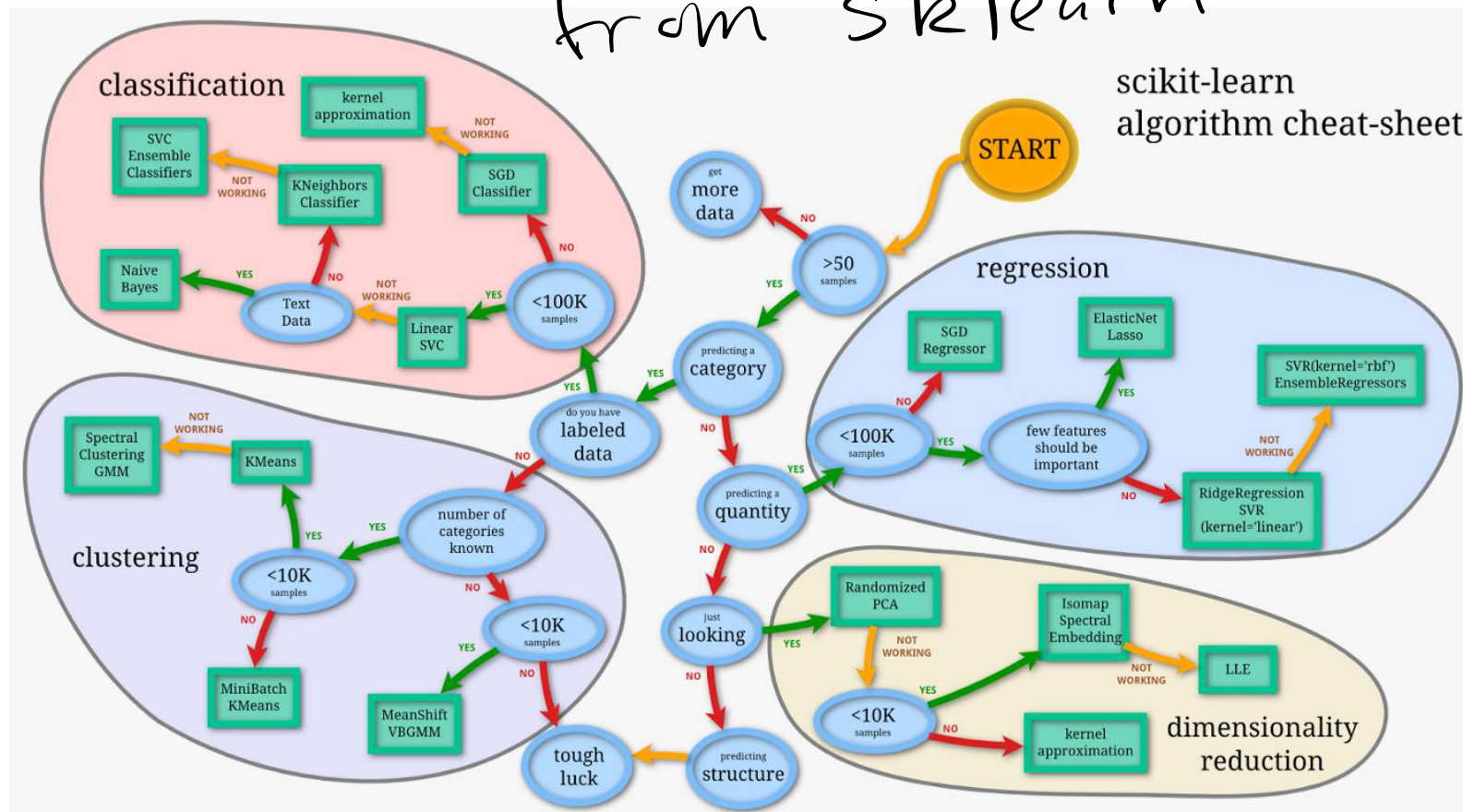
OOD \Rightarrow out of distribution
(sampling points that are distinct
from training data)

ML models \Rightarrow designed for IID
 \Rightarrow OOD sometimes correlated
w/ IID

Which model to use??

from sklearn

scikit-learn
algorithm cheat-sheet



Basics of linear regression

$$\hat{y} = f(\vec{X}) = \vec{w} \cdot \vec{X} = w_0 \vec{X}_0 + w_1 \vec{X}_1 + w_2 \vec{X}_2$$

Training \rightarrow find \vec{w} that minimize loss, $\mathcal{L}(y, \hat{y})$

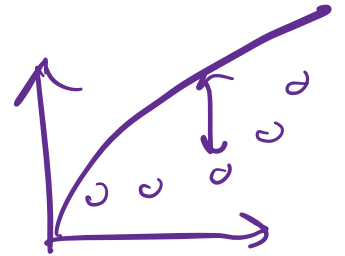
Loss functions

$$\text{OLS (MSE)} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \vec{w} \vec{X}_i)^2$$

$$\vec{w} : \min_{\vec{w}} \|\vec{X} \vec{w} - y\|_2^2$$

OLS \Rightarrow can be solved analytically
(linear algebra)

loss = error



Basics of linear regression

Regularized regression

$$\text{LASSO (L}_1\text{)} = \text{MSE} + \alpha \|\vec{w}\|_1 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \alpha \sum_{j=1}^K |w_j|$$

$K \Rightarrow \# \text{ features}$

$\alpha \Rightarrow \text{regularization parameter}$

$$\text{Ridge (L}_2\text{)} = \text{MSE} + \alpha \|\vec{w}\|_2^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 + \alpha \sum_{j=1}^K w_j^2$$

Discourages overfitting,

by penalizing large weights, w_j

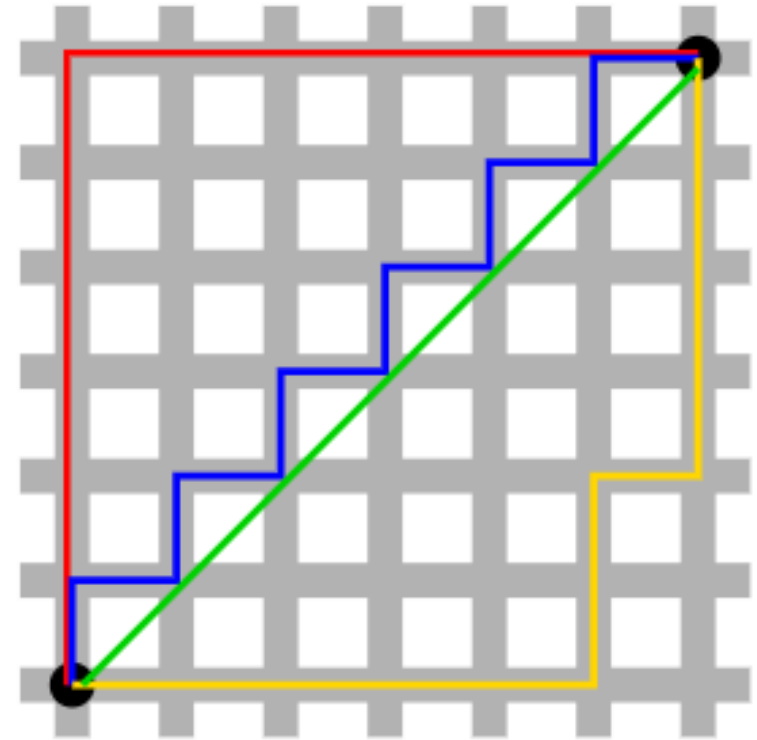
Can't be solved analytically, so need training.

L1 vs L2 norm

$$\hat{y}(w, x) = w_0 + w_1x_1 + \dots + w_px_p$$

LASSO (L1 norm): $\min_w \|Xw - y\|_2^2 + \alpha \|w\|_1$

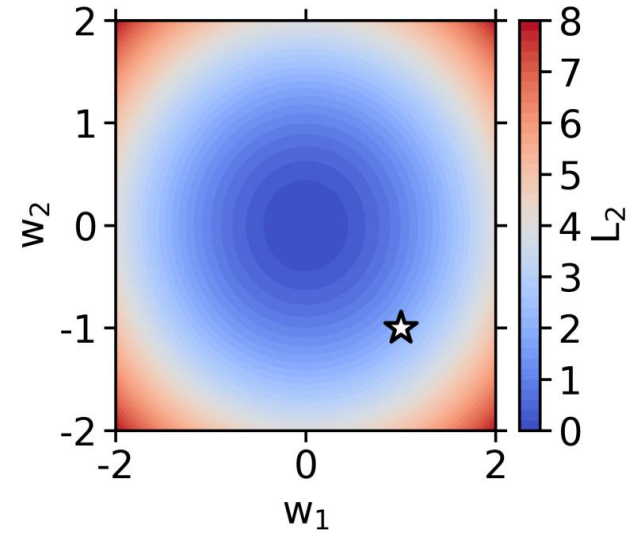
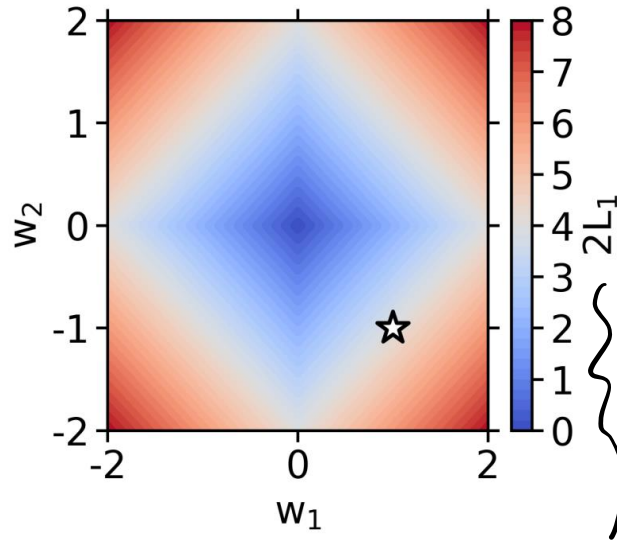
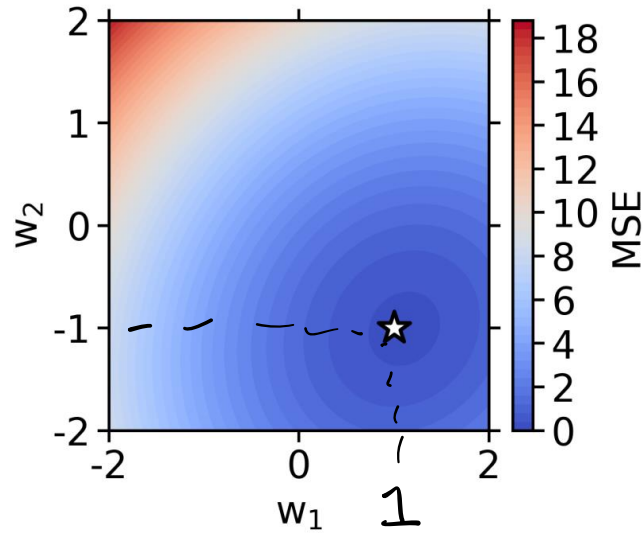
Ridge (L2 norm): $\min_w \|Xw - y\|_2^2 + \alpha \|w\|_2^2$



L1 vs L2 norm

TRUE $\rightarrow y = 1.0x_1 - 1.0x_2 + N(0, 0.25)$

$$\hat{y} = w_1x_1 + w_2x_2$$

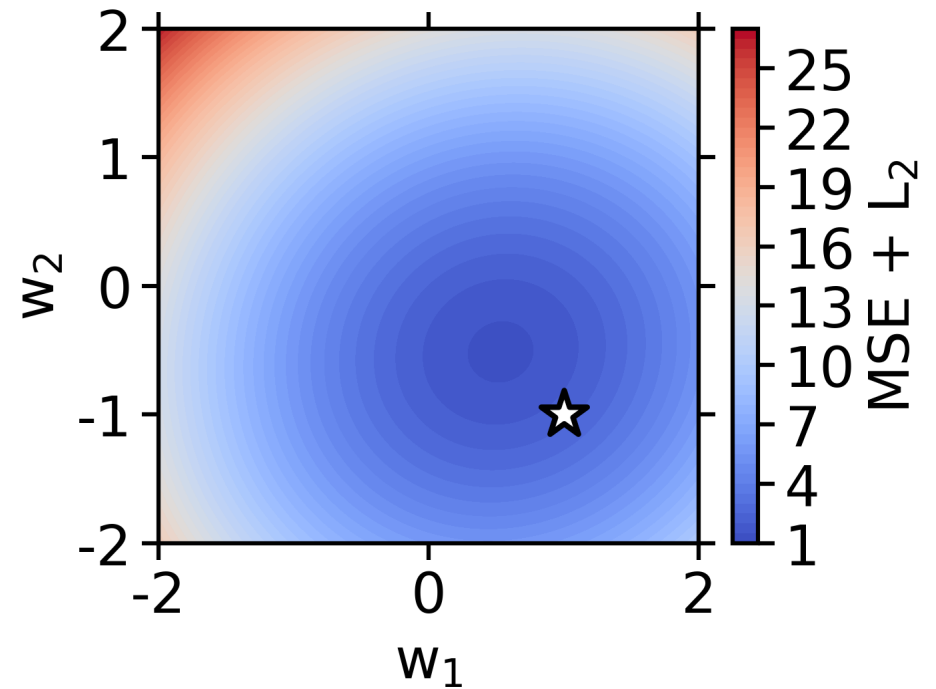
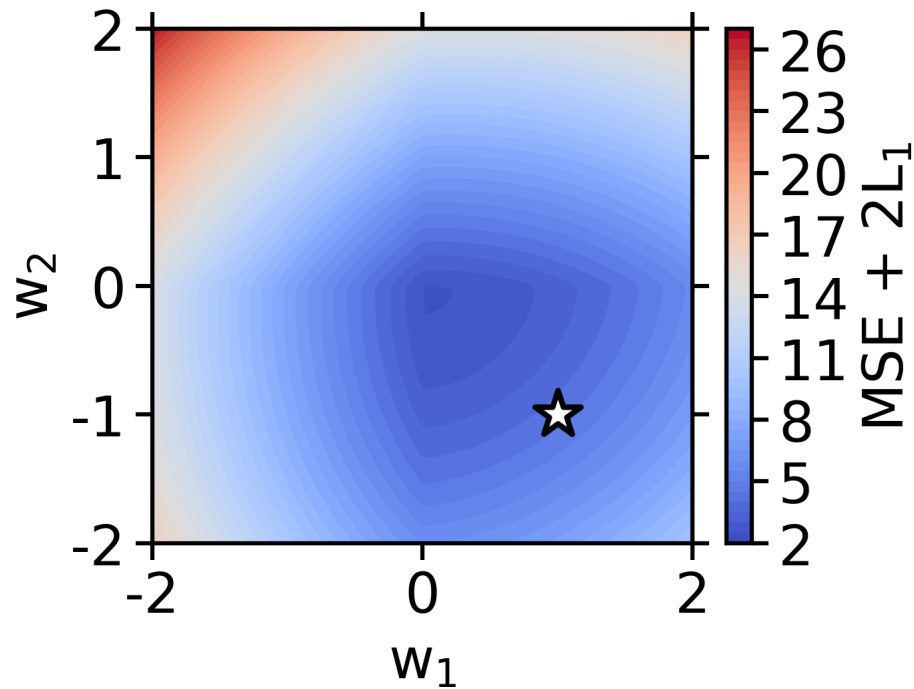


$$2L_1 = 2(|w_1| + |w_2|)$$
$$L_2 = \sqrt{w_1^2 + w_2^2}$$

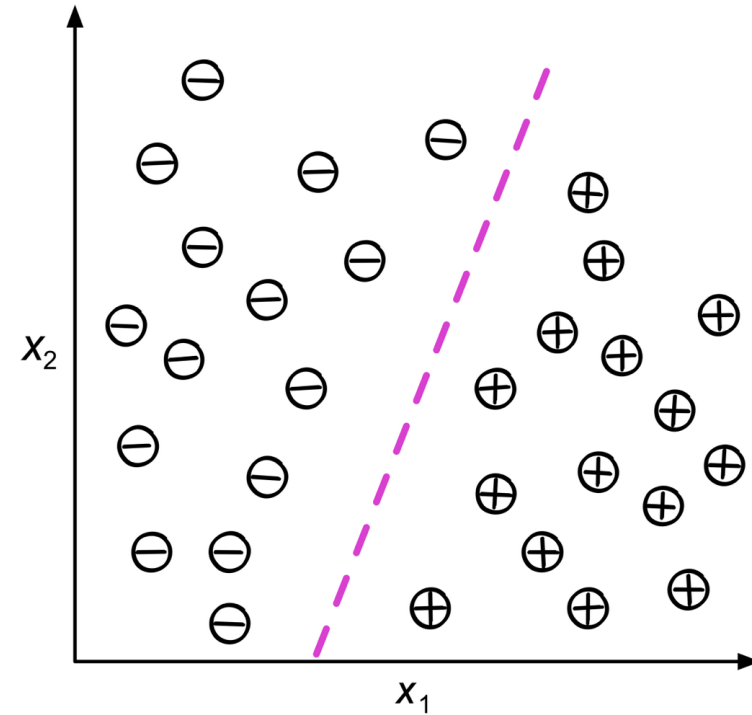
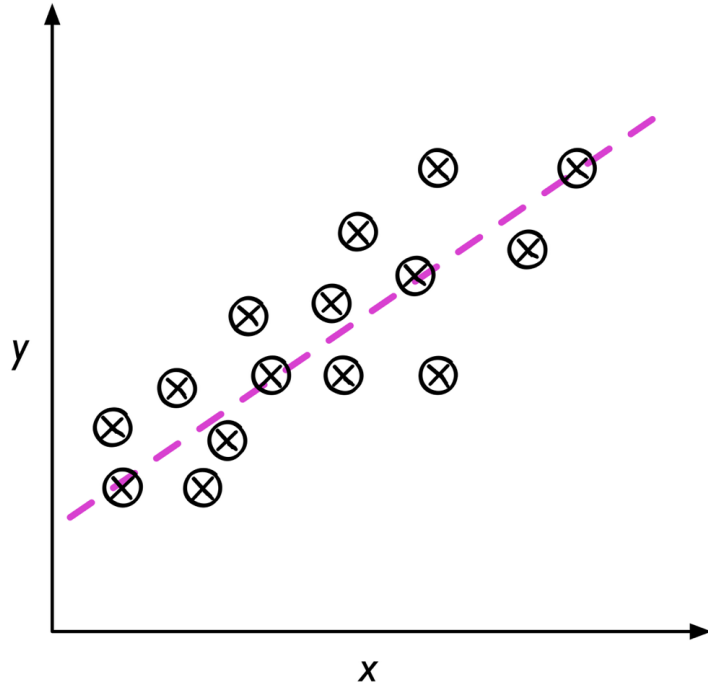
L1 vs L2 norm

$$y = 1.0x_1 - 1.0x_2 + N(0, 0.25)$$

$$\hat{y} = w_1x_1 + w_2x_2$$



How does this look for classification?



How does this look for classification?

For binary classification, y is no longer continuous, but binomial:

$$\mathbf{y} = [1, 1, 1, -1, -1, 1, -1, -1, \dots]$$

