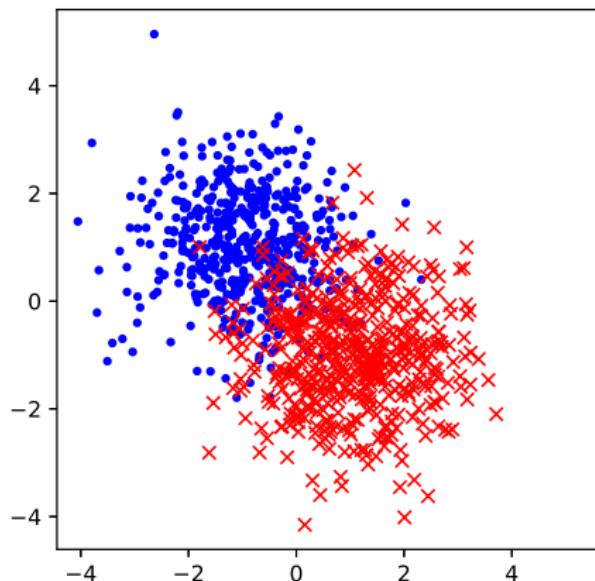


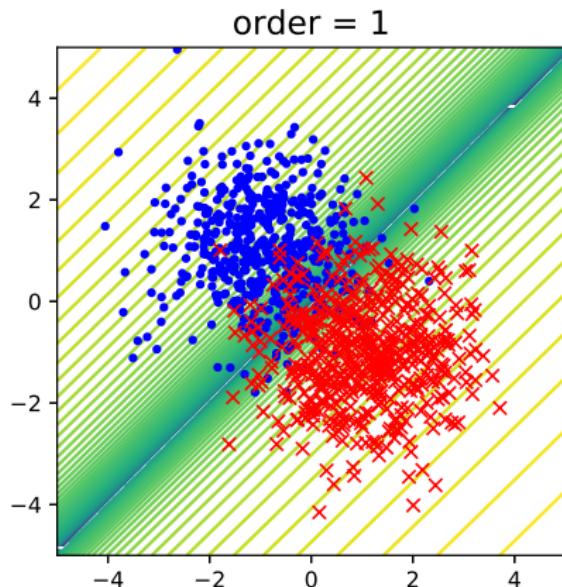
14. Generalization

- Overfitting
- Variance, bias, generalization
- Cross validation

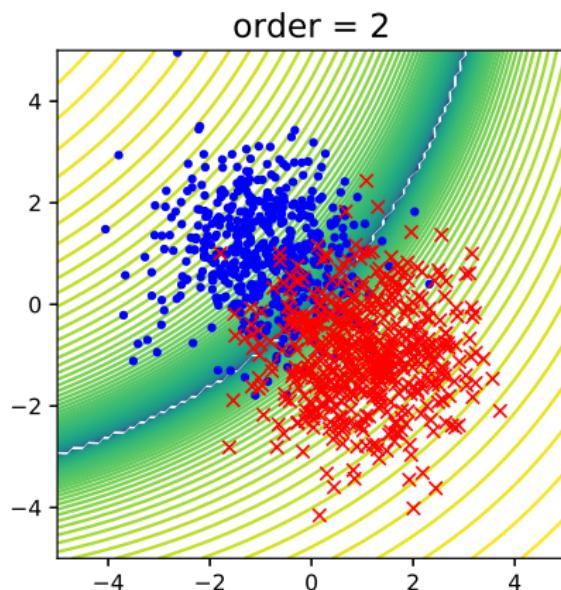
Linear classifier, polynomial basis functions



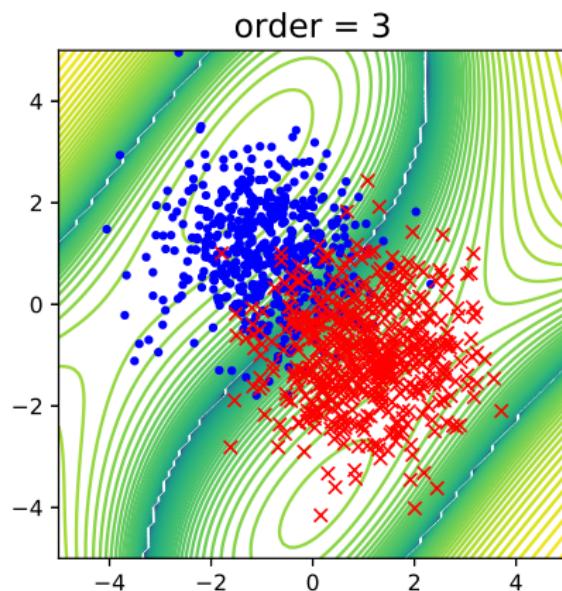
Linear classifier, polynomial basis functions



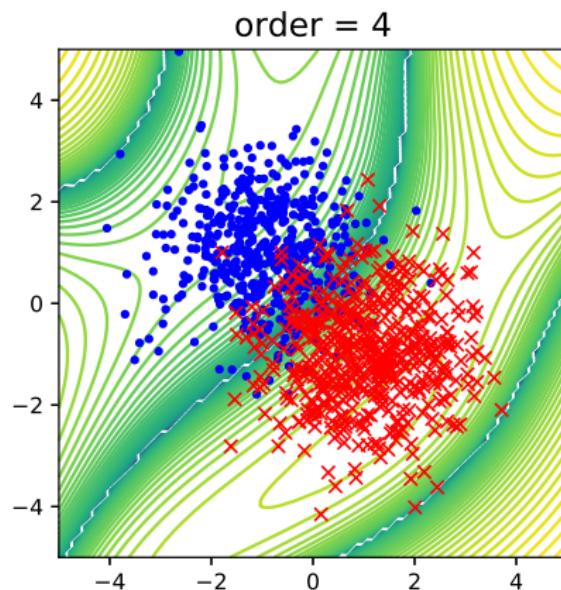
Linear classifier, polynomial basis functions



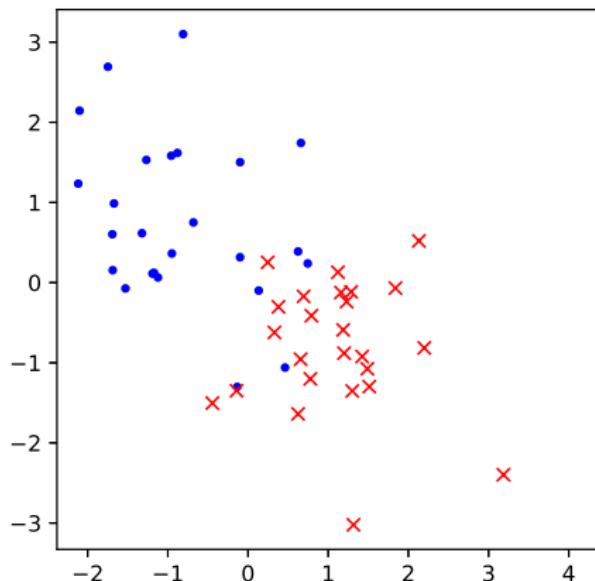
Linear classifier, polynomial basis functions



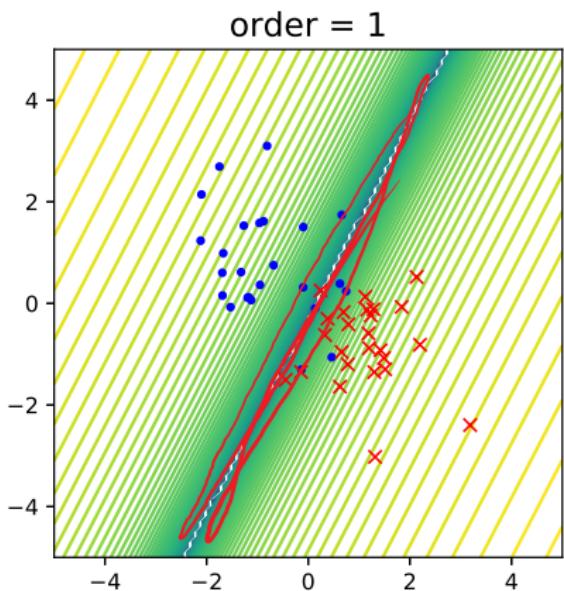
Linear classifier, polynomial basis functions



Linear classifier, polynomial basis functions

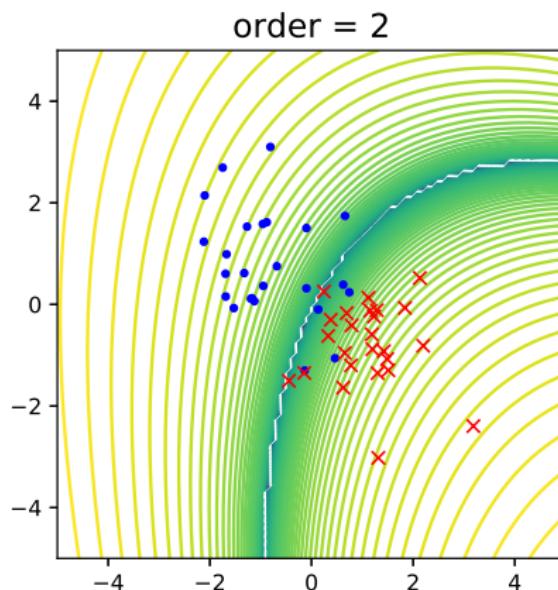


Linear classifier, polynomial basis functions

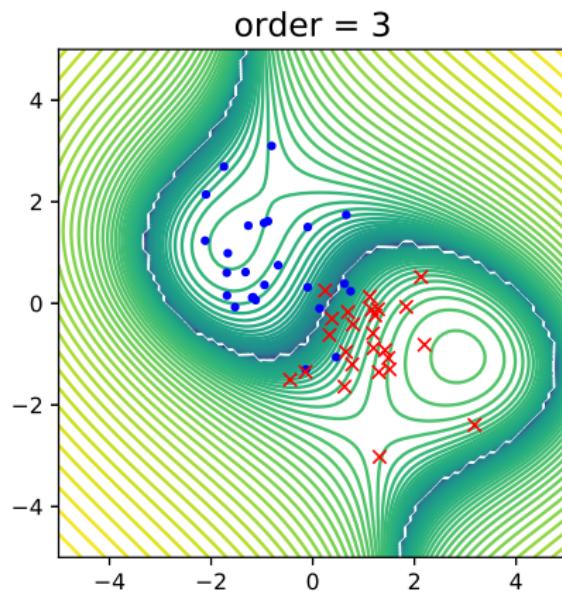


$$E[\hat{y}(x) - \hat{y}_0(x)]$$

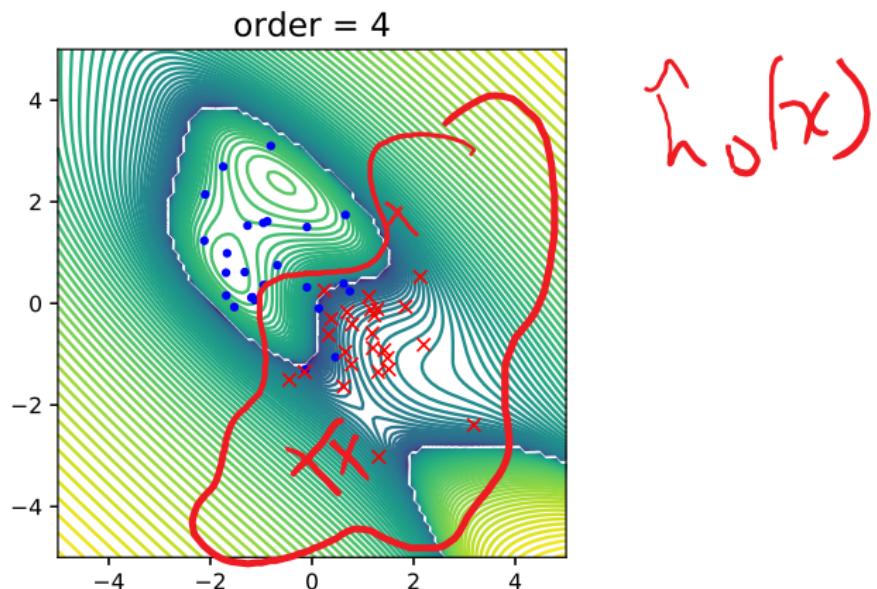
Linear classifier, polynomial basis functions



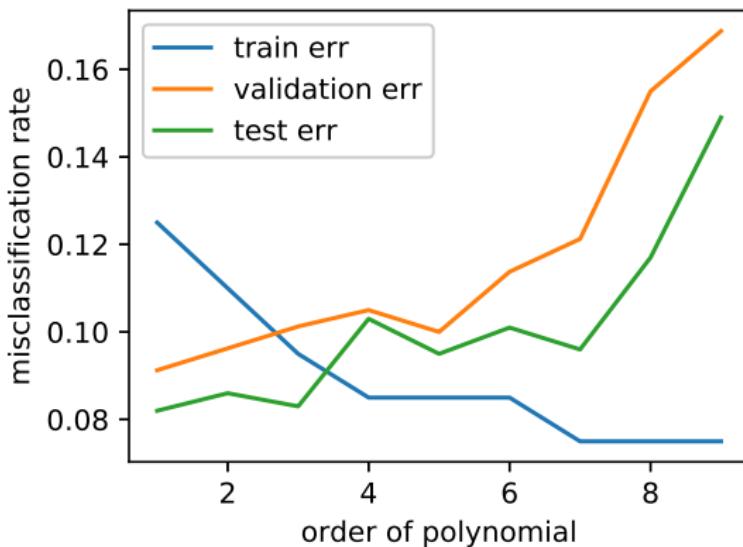
Linear classifier, polynomial basis functions



Linear classifier, polynomial basis functions

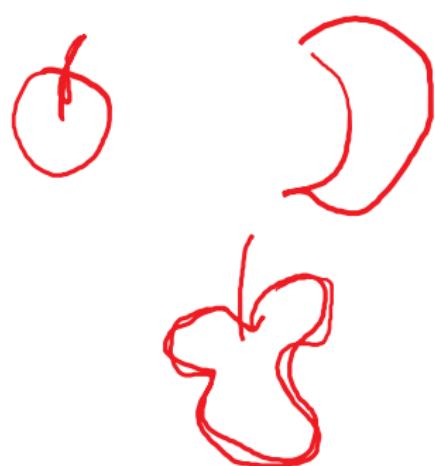


Overfitting



"perfect learn" What is overfitting?

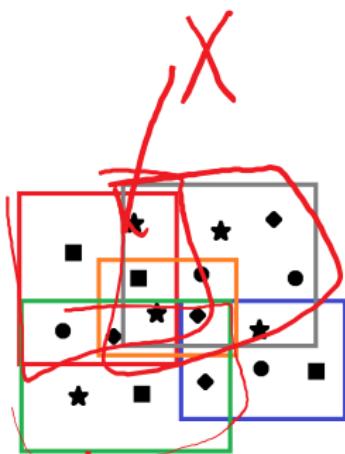
- train ↓, test doesn't
- model has "learned" noise



~~high variance~~

"out of training
distribution
error"

Model as random variable



$$\underline{D} = \{x_1, x_2, \dots, x_{|D|}\}$$

sample data
sample data
sample data
sample data
sample data



model $\{D_1\}$
model $\{D_2\}$
model $\{D_3\}$
model $\{D_4\}$
model $\{D_5\}$

Each data sample can lead to a different model

$\rightarrow D_1, D_2, \dots$

Training error vs generalization error

- Empirical risk

$$\mathcal{E}_{\mathcal{D}} := \frac{1}{|\mathcal{D}|} \sum_{x \in \mathcal{D}} \mathcal{L}(\text{pred}_{\mathcal{D}}(x), \text{label}(x))$$

random variable

"sample mean"

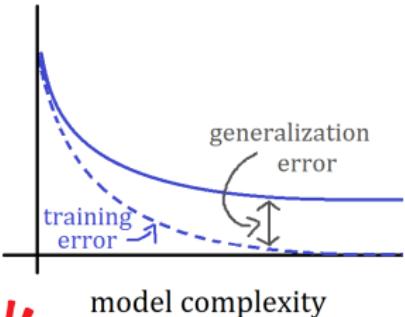
- Expected risk

$$\mathcal{E}^* := \mathbb{E}_x [\mathcal{L}(\text{pred}_{\mathcal{D}}(x), \text{label}(x))]$$

"true expectation"

- Then

$$\mathcal{E}^* = \underbrace{\mathcal{E}_{\mathcal{D}}}_{\text{training error}} + \underbrace{\mathcal{E}^* - \mathcal{E}_{\mathcal{D}}}_{\text{generalization error}}$$



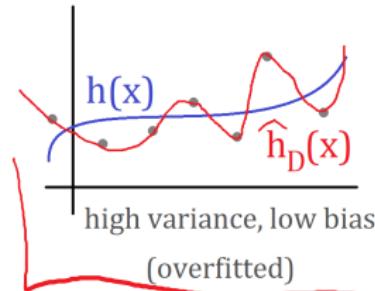
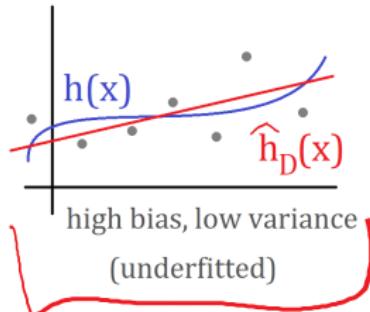
(duh!)

Bias Variance tradeoff

- True response h , trained model $\hat{h}_{\mathcal{D}}$
- bias vs variance
- Mean squared error of model on new sample x :

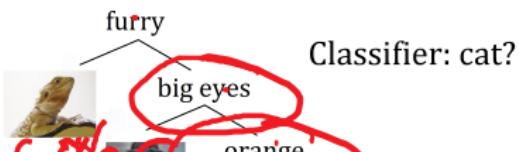
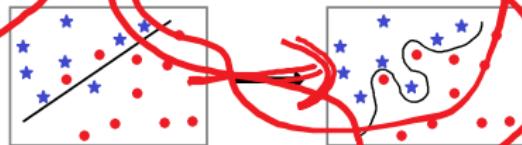
$$\mathbb{E}_{\mathcal{D}}[(h(x) - \hat{h}_{\mathcal{D}}(x))^2] =$$
$$(h(x) - \mathbb{E}_{\mathcal{D}}[\hat{h}_{\mathcal{D}}(x)])^2 + \mathbb{E}_{\mathcal{D}}(\hat{h}_{\mathcal{D}}(x)^2) - (\mathbb{E}_{\mathcal{D}}[\hat{h}_{\mathcal{D}}(x)])^2$$

$\underbrace{(h(x) - \mathbb{E}_{\mathcal{D}}[\hat{h}_{\mathcal{D}}(x)])^2}_{\text{bias}^2} + \underbrace{\mathbb{E}_{\mathcal{D}}(\hat{h}_{\mathcal{D}}(x)^2) - (\mathbb{E}_{\mathcal{D}}[\hat{h}_{\mathcal{D}}(x)])^2}_{\text{variance}}$



Example: Decision trees

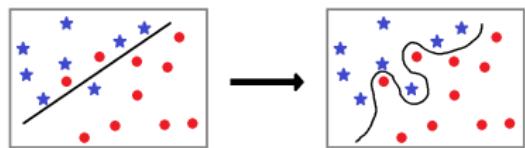
- model complexity?
- depth of tree? ←
- # samples \ll # features?
 - Can linear models overfit?



cat
black

Example: Decision trees

- model complexity?
- depth of tree?
- # samples \ll # features?
 - Can linear models overfit?
Yes! if there are far more features than samples!



Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Sample estimate:

$$\hat{\theta} = \frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} \theta$$



Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Sample estimate:

$$\hat{\theta} = \frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} \theta$$



Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Sample estimate:

$$\hat{\theta} = \frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} \theta$$



Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Sample estimate:

$$\hat{\theta} = \frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} \theta$$



Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Sample estimate:

$$\hat{\theta} = \frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} \theta$$

← *based*

$$\mathbb{E}[\hat{\theta}] = \frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} \mathbb{E}[\theta]$$



Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Sample estimate:

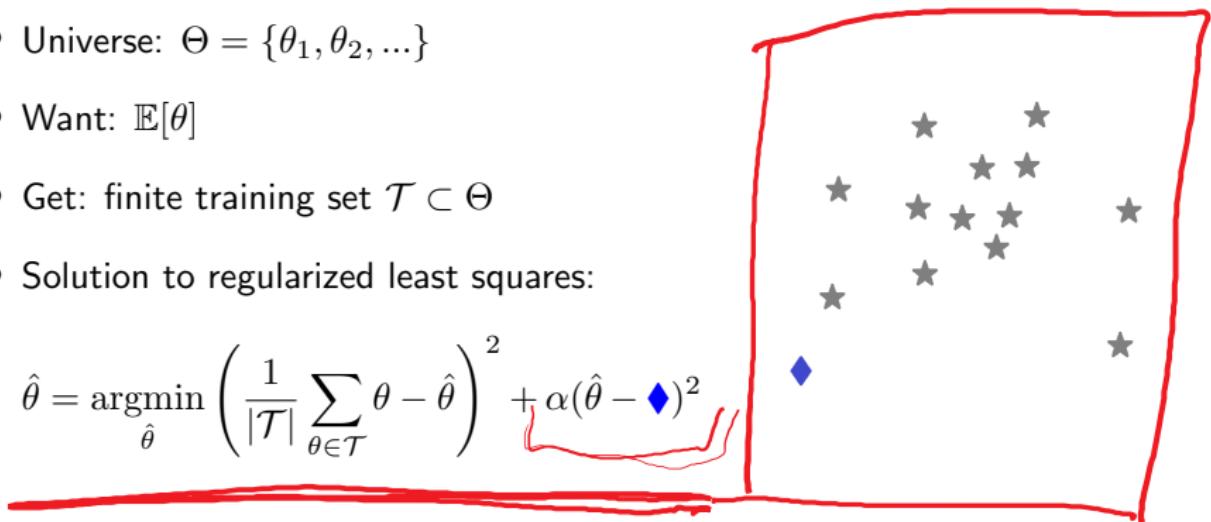
$$\hat{\theta} = \frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} \theta$$



Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Solution to regularized least squares:

$$\hat{\theta} = \operatorname{argmin}_{\hat{\theta}} \left(\frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} (\theta - \hat{\theta})^2 + \alpha(\hat{\theta} - \textcolor{blue}{\diamond})^2 \right)$$



Higher bias, lower variance

Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Solution to regularized least squares:

$$\hat{\theta} = \operatorname{argmin}_{\hat{\theta}} \left(\frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} \theta - \hat{\theta} \right)^2 + \alpha(\hat{\theta} - \diamond)^2$$



Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Solution to regularized least squares:

$$\hat{\theta} = \operatorname{argmin}_{\hat{\theta}} \left(\frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} (\theta - \hat{\theta})^2 + \alpha(\hat{\theta} - \blacklozenge)^2 \right)$$



Higher bias, lower variance

Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Solution to regularized least squares:

$$\hat{\theta} = \operatorname{argmin}_{\hat{\theta}} \left(\frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} (\theta - \hat{\theta})^2 + \alpha(\hat{\theta} - \textcolor{blue}{\diamond})^2 \right)$$



Higher bias, lower variance

Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Solution to regularized least squares:

$$\hat{\theta} = \operatorname{argmin}_{\hat{\theta}} \left(\frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} (\theta - \hat{\theta})^2 + \alpha(\hat{\theta} - \blacklozenge)^2 \right)$$

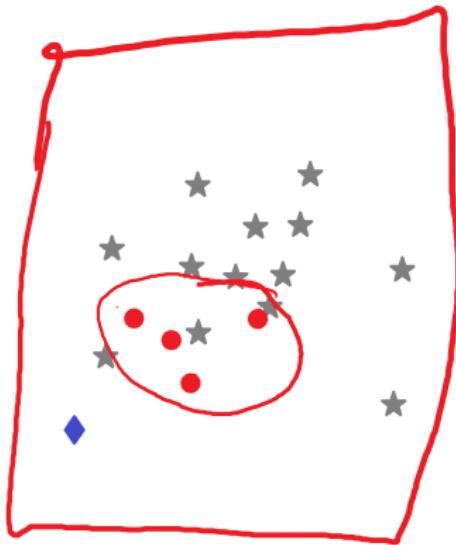


Higher bias, lower variance

Example revisited: Point estimation

- Universe: $\Theta = \{\theta_1, \theta_2, \dots\}$
- Want: $\mathbb{E}[\theta]$
- Get: finite training set $\mathcal{T} \subset \Theta$
- Solution to regularized least squares:

$$\hat{\theta} = \operatorname{argmin}_{\hat{\theta}} \left(\frac{1}{|\mathcal{T}|} \sum_{\theta \in \mathcal{T}} (\theta - \hat{\theta})^2 + \alpha(\hat{\theta} - \textcolor{blue}{\diamond})^2 \right)$$



Higher bias, lower variance

Ways to combat: Get more data!

- This is the best solution, but usually impossible/difficult in practice
- Generalization error: Assuming \mathcal{D} sampled uniformly

$$\mathcal{L}_{\mathcal{D}} \xrightarrow{|\mathcal{D}| \rightarrow +\infty} \mathcal{L}^*$$

- Variance:

$$\mathbb{E}_{\mathcal{D}}[\hat{h}_{\mathcal{D}}(x)^2] - \left(\mathbb{E}_{\mathcal{D}}[\hat{h}_{\mathcal{D}}(x)]\right)^2 = \left(\mathbb{E}_{\mathcal{D}} \underbrace{\left[\hat{h}_{\mathcal{D}}(x) - \mathbb{E}_{\mathcal{D}}[\hat{h}_{\mathcal{D}}(x)] \right]}_{\rightarrow 0 \text{ as } |\mathcal{D}| \rightarrow +\infty} \right)^2$$

- Predictor approaches best possible answer given model complexity
- In practice: increase model complexity with more data, to drive bias $\rightarrow 0$

Ways to combat: Regularization

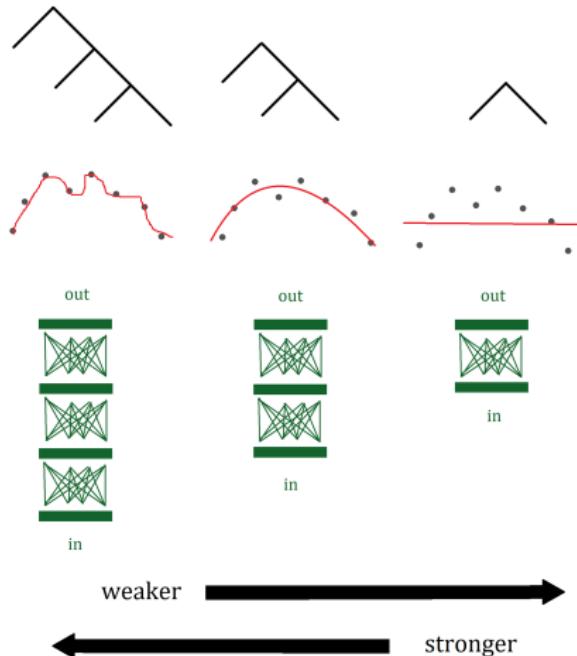
- Goal: reduce $\mathbb{E}_{x,y} \mathcal{L}_\theta(\hat{f}(x), y)$
- In practice, minimize regularized empirical risk

$$\underset{\theta}{\text{minimize}} \quad \frac{1}{m} \sum_{i=1}^m \mathcal{L}_\theta(f(x_i), y_i) + \rho \mathcal{R}(\theta)$$

- $\mathcal{R}(\theta) = \|\theta\|_2^2$ draws θ closer to 0 and improves conditioning
- $\mathcal{R}(\theta) = \|\theta\|_1$ promotes sparse θ



Ways to combat: Weakening the learner



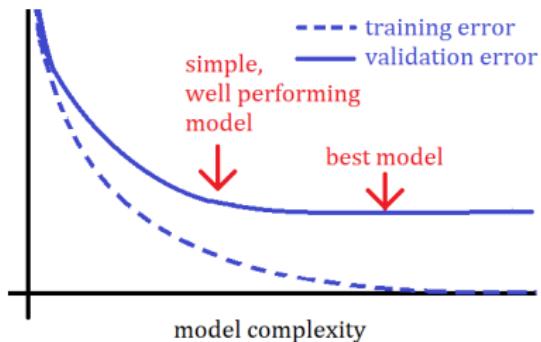
- Weaker model = larger bias, smaller variance
- Hyperparameters
 - depth of tree
 - width of neural network
 - order of polynomial
 - ...

Cross validation: Picking hyperparameters

- ① Split training data $\mathcal{D} = \mathcal{D}_T \cup \mathcal{D}_V$
- ② Sweep hyperparameters, train each model on \mathcal{D}_T
- ③ Evaluate each model on \mathcal{D}_V , pick simplest model with best validation score

Cross validation: picking hyperparameters

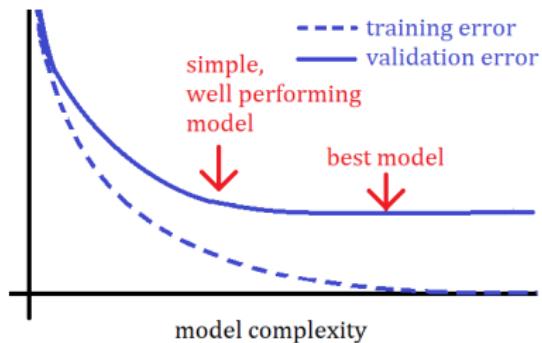
- ① Split training data $\mathcal{D} = \mathcal{D}_T \cup \mathcal{D}_V$
 - ② Sweep hyperparameters, train each model on \mathcal{D}_T
 - ③ Evaluate each model on \mathcal{D}_V , pick simplest model with best validation score
- “K-fold cross validation, 75 /25 split” = repeat 1-3 K times with $\mathcal{D} = \underbrace{\mathcal{D}_T}_{75\%} \cup \underbrace{\mathcal{D}_V}_{25\%}$, randomly partitioned
 - Key: never use test set in cross validation! (why?)



Cross validation: picking hyperparameters

- ① Split training data $\mathcal{D} = \mathcal{D}_T \cup \mathcal{D}_V$
- ② Sweep hyperparameters, train each model on \mathcal{D}_T
- ③ Evaluate each model on \mathcal{D}_V , pick simplest model with best validation score

- “K-fold cross validation, 75 /25 split” =
repeat 1-3 K times with $\mathcal{D} = \underbrace{\mathcal{D}_T}_{75\%} \cup \underbrace{\mathcal{D}_V}_{25\%}$, randomly partitioned
- Key: never use test set in cross validation! (why?)
Hyperparameters too can be overfitted!



K-fold cross validation

- Divide the data into K disjoint subsets
 - Train on union of $K - 1$ subsets
 - Test on the left-out set
 - Repeat K times
- Leave one out CV (LOOCV) is K -fold CV with $K = \#$ of data points

Summary

- Training error: how well your model fit data you saw
- Generalization error: how consistent your model behaves on unseen data
- Balancing act

Total error = training error + generalization error

- Ways to improve generalization error
 - More data
 - Weakening the model
 - Cross validation