

COMS 4771 Lecture 12

1. Introduction to learning theory
2. Cross validation

INTRODUCTION TO LEARNING THEORY

Basic setting

- ▶ Training data S is an iid sample from some fixed but unknown probability distribution P over space of labeled examples $\mathcal{X} \times \mathcal{Y}$.

Basic setting

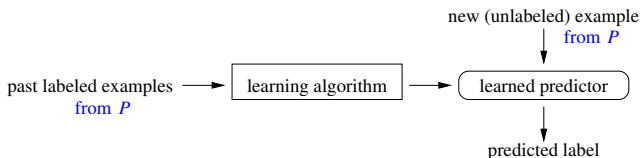
- ▶ Training data S is an iid sample from some fixed but unknown probability distribution P over space of labeled examples $\mathcal{X} \times \mathcal{Y}$.
- ▶ A learning algorithm takes S as input and returns a predictor $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$.

Basic setting

- ▶ Training data S is an iid sample from some fixed but unknown probability distribution P over space of labeled examples $\mathcal{X} \times \mathcal{Y}$.
- ▶ A learning algorithm takes S as input and returns a predictor $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$.
- ▶ Benchmark: (true) prediction error $\text{err}(\hat{f})$.
(Recall, $\text{err}(f) = \Pr(f(X) \neq Y)$, where (X, Y) has distribution P .)

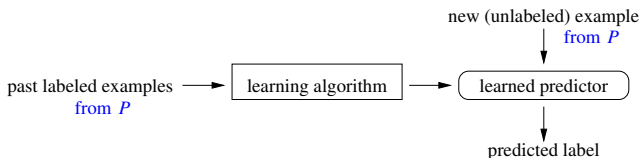
Basic setting

- ▶ Training data S is an iid sample from some fixed but unknown probability distribution P over space of labeled examples $\mathcal{X} \times \mathcal{Y}$.
- ▶ A learning algorithm takes S as input and returns a predictor $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$.
- ▶ Benchmark: (true) prediction error $\text{err}(\hat{f})$.
(Recall, $\text{err}(f) = \Pr(f(X) \neq Y)$, where (X, Y) has distribution P .)



Basic setting

- ▶ Training data S is an iid sample from some fixed but unknown probability distribution P over space of labeled examples $\mathcal{X} \times \mathcal{Y}$.
- ▶ A learning algorithm takes S as input and returns a predictor $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$.
- ▶ Benchmark: (true) prediction error $\text{err}(\hat{f})$.
(Recall, $\text{err}(f) = \Pr(f(X) \neq Y)$, where (X, Y) has distribution P .)

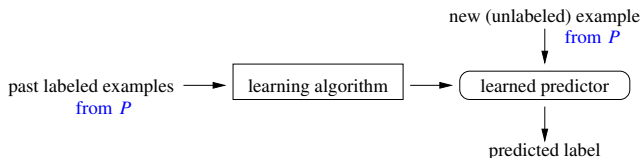


In this setting, can any learning algorithm always provide a non-trivial guarantee on the error of the predictor it returns?

STATISTICAL LEARNING

Basic setting

- ▶ Training data S is an iid sample from some fixed but unknown probability distribution P over space of labeled examples $\mathcal{X} \times \mathcal{Y}$.
- ▶ A learning algorithm takes S as input and returns a predictor $\hat{f}: \mathcal{X} \rightarrow \mathcal{Y}$.
- ▶ Benchmark: (true) prediction error $\text{err}(\hat{f})$.
(Recall, $\text{err}(f) = \Pr(f(X) \neq Y)$, where (X, Y) has distribution P .)



In this setting, can any learning algorithm always provide a non-trivial guarantee on the error of the predictor it returns?

No: some assumptions/conditions are required.
(“No free lunch” theorem)

ASSUMPTION: REALIZABILITY

Consider binary classification with $\mathcal{Y} = \{0, 1\}$.

Realizability assumption: Assume that, for a given class \mathcal{F} of functions from $\mathcal{X} \rightarrow \{0, 1\}$, there exists $f^* \in \mathcal{F}$ such that $\text{err}(f^*) = 0$.

ASSUMPTION: REALIZABILITY

Consider binary classification with $\mathcal{Y} = \{0, 1\}$.

Realizability assumption: Assume that, for a given class \mathcal{F} of functions from $\mathcal{X} \rightarrow \{0, 1\}$, there exists $f^* \in \mathcal{F}$ such that $\text{err}(f^*) = 0$.

(This implies that f^* is the *Bayes classifier*.)

ASSUMPTION: REALIZABILITY

Consider binary classification with $\mathcal{Y} = \{0, 1\}$.

Realizability assumption: Assume that, for a given class \mathcal{F} of functions from $\mathcal{X} \rightarrow \{0, 1\}$, there exists $f^* \in \mathcal{F}$ such that $\text{err}(f^*) = 0$.

(This implies that f^* is the *Bayes classifier*.)

Examples of function classes \mathcal{F} :

- ▶ Rectangles in $\mathcal{X} = \mathbb{R}^2$:

$$f_{((a,b),(c,d))}(\mathbf{x}) = \mathbb{1}\{a \leq x_1 \leq b \text{ and } c \leq x_2 \leq d\}.$$

- ▶ Monotone conjunctions in $\mathcal{X} = \{0, 1\}^d$:

$$f_V = \mathbb{1}\{x_i = 1 \text{ for all } i \in V\}.$$

- ▶ Linear classifiers in $\mathcal{X} = \mathbb{R}^d$.

ASSUMPTION: REALIZABILITY

Consider binary classification with $\mathcal{Y} = \{0, 1\}$.

Realizability assumption: Assume that, for a given class \mathcal{F} of functions from $\mathcal{X} \rightarrow \{0, 1\}$, there exists $f^* \in \mathcal{F}$ such that $\text{err}(f^*) = 0$.

(This implies that f^* is the *Bayes classifier*.)

Examples of function classes \mathcal{F} :

- Rectangles in $\mathcal{X} = \mathbb{R}^2$:

$$f_{((a,b),(c,d))}(\mathbf{x}) = \mathbb{1}\{a \leq x_1 \leq b \text{ and } c \leq x_2 \leq d\}.$$

- Monotone conjunctions in $\mathcal{X} = \{0, 1\}^d$:

$$f_V = \mathbb{1}\{x_i = 1 \text{ for all } i \in V\}.$$

- Linear classifiers in $\mathcal{X} = \mathbb{R}^d$.

Realizable setting is essentially the setup of **PAC Learning**, a theoretical model of learning introduced by L. Valiant (1984).

LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

Pick a consistent classifier (a special case of ERM):

Given training data S , return any $f \in \mathcal{F}$ such that $\text{err}(f, S) = 0$.

Always possible under realizability assumption!

LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

Pick a consistent classifier (a special case of ERM):

Given training data S , return any $\hat{f} \in \mathcal{F}$ such that $\text{err}(\hat{f}, S) = 0$.

Always possible under realizability assumption!

Example: rectangles in \mathbb{R}^2 .

LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

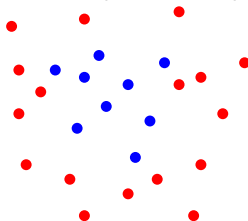
Pick a consistent classifier (a special case of ERM):

Given training data S , return any $\hat{f} \in \mathcal{F}$ such that $\text{err}(\hat{f}, S) = 0$.

Always possible under realizability assumption!

Example: rectangles in \mathbb{R}^2 .

(Blue dots = positive examples.)



LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

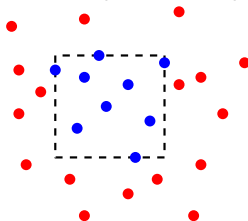
Pick a consistent classifier (a special case of ERM):

Given training data S , return any $f \in \mathcal{F}$ such that $\text{err}(f, S) = 0$.

Always possible under realizability assumption!

Example: rectangles in \mathbb{R}^2 .

(Blue dots = positive examples.)



LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

Pick a consistent classifier (a special case of ERM):

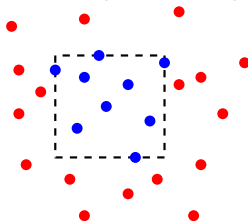
Given training data S , return any $\hat{f} \in \mathcal{F}$ such that $\text{err}(\hat{f}, S) = 0$.

Always possible under realizability assumption!

Example: rectangles in \mathbb{R}^2 .

Example: monotone conjunctions in $\{0, 1\}^d$.

(Blue dots = positive examples.)



LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

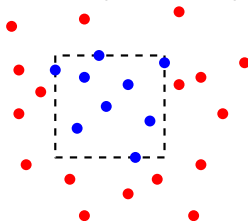
Pick a consistent classifier (a special case of ERM):

Given training data S , return any $\hat{f} \in \mathcal{F}$ such that $\text{err}(\hat{f}, S) = 0$.

Always possible under realizability assumption!

Example: rectangles in \mathbb{R}^2 .

(Blue dots = positive examples.)



Example: monotone conjunctions in $\{0, 1\}^d$.

Return monotone conjunction f_V with

$$V := [d] \setminus \left(\bigcup_{(\mathbf{x}, +1) \in S} \{i \in [d] : x_i = 0\} \right).$$

LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

Pick a consistent classifier (a special case of ERM):

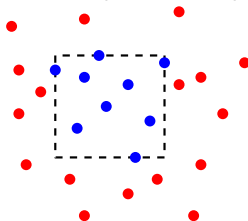
Given training data S , return any $\hat{f} \in \mathcal{F}$ such that $\text{err}(\hat{f}, S) = 0$.

since it's
reliable!

Always possible under realizability assumption!

Example: rectangles in \mathbb{R}^2 .

(Blue dots = positive examples.)



Example: monotone conjunctions in $\{0, 1\}^d$.

Return monotone conjunction f_V with

$$V := [d] \setminus \left(\bigcup_{(\mathbf{x}, +1) \in S} \{i \in [d] : x_i = 0\} \right).$$

- ▶ Start with $V := [d]$.
- ▶ For each positive example $(\mathbf{x}, +1) \in S$, remove all $i \in [d]$ from V s.t. $x_i = 0$.

LEARNING IN THE REALIZABLE SETTING

What is a sensible learning algorithm for the realizable setting?

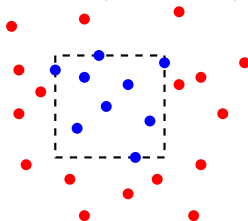
Pick a consistent classifier (a special case of ERM):

Given training data S , return any $f \in \mathcal{F}$ such that $\text{err}(\hat{f}, S) = 0$.

Always possible under realizability assumption!

Example: rectangles in \mathbb{R}^2 .

(Blue dots = positive examples.)



Example: monotone conjunctions in $\{0, 1\}^d$.

Return monotone conjunction f_V with

$$V := [d] \setminus \left(\bigcup_{(\mathbf{x}, +1) \in S} \{i \in [d] : x_i = 0\} \right).$$

- ▶ Start with $V := [d]$.
- ▶ For each positive example $(\mathbf{x}, +1) \in S$, remove all $i \in [d]$ from V s.t. $x_i = 0$.

Example: linear classifiers in $\mathbb{R}^d \rightarrow$ linear programming, Perceptron, or SVM.

CONSISTENT CLASSIFIER ALGORITHM

Is the “Consistent Classifier Algorithm” any good in the realizable setting?

- ▶ Would like to show that typically, the error of the returned classifier $\hat{f} \in \mathcal{F}$ goes to zero as the number of training data increases.

CONSISTENT CLASSIFIER ALGORITHM

Is the “Consistent Classifier Algorithm” any good in the realizable setting?

- ▶ Would like to show that typically, the error of the returned classifier $\hat{f} \in \mathcal{F}$ goes to zero as the number of training data increases.
- ▶ Formally, want for any $\delta \in (0, 1)$,

$$\Pr \left[\text{err}(\hat{f}) \leq \varepsilon(|S|) \right] \geq 1 - \delta$$

for some $\varepsilon(n)$ satisfying

$$\lim_{n \rightarrow \infty} \varepsilon(n) = 0.$$

Note: $\varepsilon(n)$ may also depend on $\delta, \mathcal{F}, f^*, \dots$

CONSISTENT CLASSIFIER ALGORITHM

Is the “Consistent Classifier Algorithm” any good in the realizable setting?

- ▶ Would like to show that typically, the error of the returned classifier $\hat{f} \in \mathcal{F}$ goes to zero as the number of training data increases.
- ▶ Formally, want for any $\delta \in (0, 1)$,

$$\Pr \left[\text{err}(\hat{f}) \leq \varepsilon(|S|) \right] \geq 1 - \delta$$

for some $\varepsilon(n)$ satisfying

$$\lim_{n \rightarrow \infty} \varepsilon(n) = 0.$$

Note: $\varepsilon(n)$ may also depend on δ , \mathcal{F} , f^* , \dots

Note 2: Could also ask for

$$\mathbb{E} \left[\text{err}(\hat{f}) \right] \leq \varepsilon(|S|) \rightarrow 0$$

(i.e., *statistical consistency*).

CONSISTENT CLASSIFIER ALGORITHM

Analysis of “Consistent Classifier Algorithm”

We know that for any $f: \mathcal{X} \rightarrow \{0, 1\}$,

$$\Pr \left[\text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(1/\delta)}{|S|}} + \frac{2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

(Upper limit of confidence interval for a coin bias based on Chernoff bounds.)

CONSISTENT CLASSIFIER ALGORITHM

error rate on sample is lower than
the real error

Analysis of “Consistent Classifier Algorithm”

We know that for any $f: \mathcal{X} \rightarrow \{0, 1\}$,

$$\Pr \left[\text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(1/\delta)}{|S|}} + \frac{2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

(Upper limit of confidence interval for a coin bias based on Chernoff bounds.)

Question:

Does this apply to classifier $\hat{f} \in \mathcal{F}$ returned by the algorithm?

CONSISTENT CLASSIFIER ALGORITHM

Analysis of “Consistent Classifier Algorithm”

We know that for any $f: \mathcal{X} \rightarrow \{0, 1\}$,

$$\Pr \left[\text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(1/\delta)}{|S|}} + \frac{2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

(Upper limit of confidence interval for a coin bias based on Chernoff bounds.)

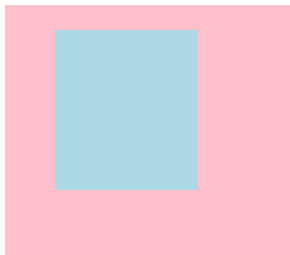
Question:

Does this apply to classifier $\hat{f} \in \mathcal{F}$ returned by the algorithm?

Generally no \hat{f} is not fixed; it's chosen based on (random) training data S .

DETOUR: OVERFITTING

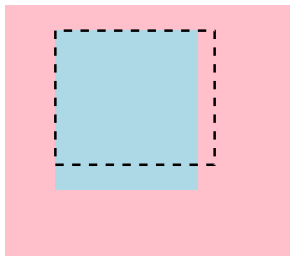
Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Data distribution P over $\mathbb{R}^2 \times \{0, 1\}$.
(blue = positive mass)

DETOUR: OVERFITTING

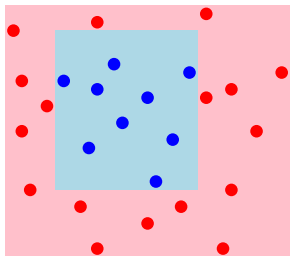
Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Particular rectangle function $f_1 \in \mathcal{F}$.

DETOUR: OVERFITTING

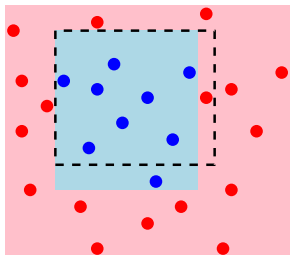
Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Random sample S from P .

DETOUR: OVERFITTING

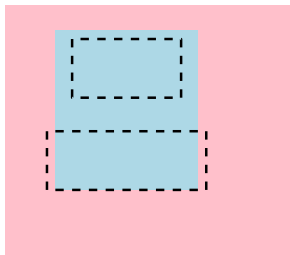
Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Particular rectangle function $f_1 \in \mathcal{F}$ and \mathcal{S} .
 $\text{err}(f_1, \mathcal{S}) \approx \text{err}(f_1)$

DETOUR: OVERFITTING

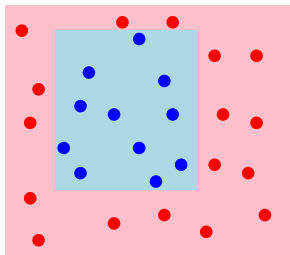
Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Union of rectangles function $f_2 \in \mathcal{F}$.

DETOUR: OVERFITTING

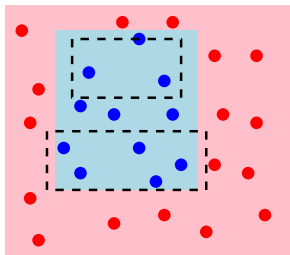
Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Random sample S' from P .

DETOUR: OVERFITTING

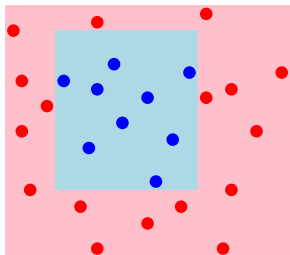
Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Union of rectangles function $f_2 \in \mathcal{F}$ on S' .
 $\text{err}(f_2, S') \approx \text{err}(f_2)$

DETOUR: OVERFITTING

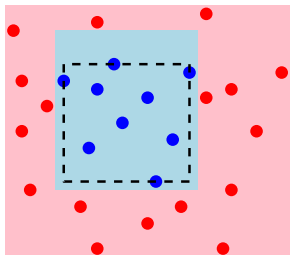
Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Back to first sample S from P .

DETOUR: OVERFITTING

Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .

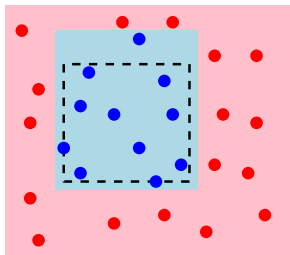


Rectangle function $\hat{f}_{3,S} \in \mathcal{F}$ on S .

$$0 = \text{err}(\hat{f}_{3,S}, S) < \text{err}(\hat{f}_{3,S})$$

DETOUR: OVERFITTING

Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .

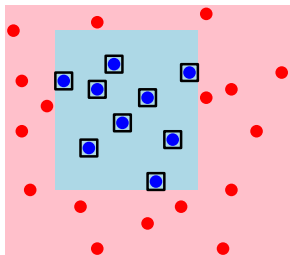


Rectangle function $\hat{f}_{3,S} \in \mathcal{F}$ on S' .

$$0 < \text{err}(\hat{f}_{3,S}, S') \approx \text{err}(\hat{f}_{3,S})$$

DETOUR: OVERFITTING

Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .

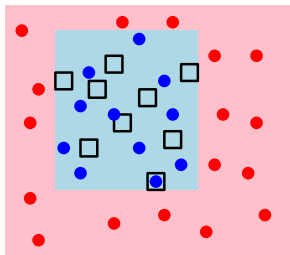


Union of rectangles $\hat{f}_{4,S} \in \mathcal{F}$ on S .

$$0 = \text{err}(\hat{f}_{4,S}, S) \ll \text{err}(\hat{f}_{4,S})$$

DETOUR: OVERFITTING

Consider $\mathcal{F} =$ union of up to 9 rectangles in \mathbb{R}^2 .



Union of rectangles $\hat{f}_{4,S} \in \mathcal{F}$ on S' .

$$0 \ll \text{err}(\hat{f}_{4,S}, S') \approx \text{err}(\hat{f}_{4,S})$$

CONSISTENT CLASSIFIER ALGORITHM

Analysis of “Consistent Classifier Algorithm”

Moral: for a classifier *chosen using the training data*, training error is not an unbiased estimate of true error.

Text

CONSISTENT CLASSIFIER ALGORITHM

Analysis of “Consistent Classifier Algorithm”

Moral: for a classifier *chosen using the training data*, **training error is not an unbiased estimate of true error.**

For all $f: \mathcal{X} \rightarrow \{0, 1\}$,

$$\Pr \left[\text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(1/\delta)}{|S|}} + \frac{2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

(Upper limit of confidence interval for a coin bias based on Chernoff bounds.)

CONSISTENT CLASSIFIER ALGORITHM

not unbiased!

Analysis of “Consistent Classifier Algorithm”

Moral: for a classifier *chosen using the training data*, **training error is not an unbiased estimate of true error**.

For all $f: \mathcal{X} \rightarrow \{0, 1\}$,

$$\Pr \left[\text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(1/\delta)}{|S|}} + \frac{2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

(Upper limit of confidence interval for a coin bias based on Chernoff bounds.)

Overkill solution: ensure upper confidence bounds hold for *all* $f \in \mathcal{F}$ *simultaneously*, with probability $\geq 1 - \delta$.

CONSISTENT CLASSIFIER ALGORITHM FOR FINITE \mathcal{F}

Analysis of Consistent Classifier Algorithm for finite \mathcal{F}

Union bound: for any countable sequence of events $\mathcal{E}_1, \mathcal{E}_2, \dots$,

$$\Pr \left[\bigcup_{i \geq 1} \mathcal{E}_i \right] \leq \sum_{i \geq 1} \Pr[\mathcal{E}_i].$$

CONSISTENT CLASSIFIER ALGORITHM FOR FINITE \mathcal{F}

Analysis of Consistent Classifier Algorithm for finite \mathcal{F}

Union bound: for any countable sequence of events $\mathcal{E}_1, \mathcal{E}_2, \dots$,

$$\Pr \left[\bigcup_{i \geq 1} \mathcal{E}_i \right] \leq \sum_{i \geq 1} \Pr[\mathcal{E}_i].$$

Apply to events \mathcal{E}_f for $f \in \mathcal{F}$ given by

$$\mathcal{E}_f := \left\{ \text{err}(f) > \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(1/\delta)}{|S|}} + \frac{2 \ln(1/\delta)}{|S|} \right\}.$$

(From last slide: $\Pr[\mathcal{E}_f] \leq \delta$ for each $f \in \mathcal{F}$.)

CONSISTENT CLASSIFIER ALGORITHM FOR FINITE \mathcal{F}

Analysis of Consistent Classifier Algorithm for finite \mathcal{F}

Union bound: for any countable sequence of events $\mathcal{E}_1, \mathcal{E}_2, \dots$,

$$\Pr \left[\bigcup_{i \geq 1} \mathcal{E}_i \right] \leq \sum_{i \geq 1} \Pr[\mathcal{E}_i].$$

Apply to events \mathcal{E}_f for $f \in \mathcal{F}$ given by

$$\mathcal{E}_f := \left\{ \text{err}(f) > \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(1/\delta)}{|S|}} + \frac{2 \ln(1/\delta)}{|S|} \right\}.$$

(From last slide: $\Pr[\mathcal{E}_f] \leq \delta$ for each $f \in \mathcal{F}$.)

Therefore, $\Pr[\bigcup_{f \in \mathcal{F}} \mathcal{E}_f] \leq |\mathcal{F}| \delta \dots$

CONSISTENT CLASSIFIER ALGORITHM FOR FINITE \mathcal{F}

Analysis of Consistent Classifier Algorithm for finite \mathcal{F}

Union bound: for any countable sequence of events $\mathcal{E}_1, \mathcal{E}_2, \dots$,

intersection!!
$$\Pr \left[\bigcup_{i \geq 1} \mathcal{E}_i \right] \leq \sum_{i \geq 1} \Pr[\mathcal{E}_i].$$

Apply to events \mathcal{E}_f for $f \in \mathcal{F}$ given by

$$\mathcal{E}_f := \left\{ \text{err}(f) > \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(1/\delta)}{|S|}} + \frac{2 \ln(1/\delta)}{|S|} \right\}.$$

(From last slide: $\Pr[\mathcal{E}_f] \leq \delta$ for each $f \in \mathcal{F}$.)

Therefore, $\Pr[\bigcup_{f \in \mathcal{F}} \mathcal{E}_f] \leq |\mathcal{F}| \delta \dots$ i.e., (replacing δ with $\delta/|\mathcal{F}|$)

$$\Pr \left[\forall f \in \mathcal{F} \bullet \text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(|\mathcal{F}|/\delta)}{|S|}} + \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

CONSISTENT CLASSIFIER ALGORITHM FOR FINITE \mathcal{F}

From last slide:

$$\Pr \left[\forall f \in \mathcal{F} \bullet \text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(|\mathcal{F}|/\delta)}{|S|}} + \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

CONSISTENT CLASSIFIER ALGORITHM FOR FINITE \mathcal{F}

From last slide:

$$\Pr \left[\forall f \in \mathcal{F} \bullet \text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(|\mathcal{F}|/\delta)}{|S|}} + \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

Since the Consistent Classifier Algorithm returns $\hat{f} \in \mathcal{F}$, we know that

$$\Pr \left[\text{err}(\hat{f}) \leq \text{err}(\hat{f}, S) + \sqrt{\frac{2 \text{err}(\hat{f}, S) \ln(|\mathcal{F}|/\delta)}{|S|}} + \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

CONSISTENT CLASSIFIER ALGORITHM FOR FINITE \mathcal{F}

From last slide:

$$\Pr \left[\forall f \in \mathcal{F} \bullet \text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(|\mathcal{F}|/\delta)}{|S|}} + \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

Since the Consistent Classifier Algorithm returns $\hat{f} \in \mathcal{F}$, we know that

$$\Pr \left[\text{err}(\hat{f}) \leq \text{err}(\hat{f}, S) + \sqrt{\frac{2 \text{err}(\hat{f}, S) \ln(|\mathcal{F}|/\delta)}{|S|}} + \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

key!!!

By definition of \hat{f} , $\text{err}(\hat{f}, S) = 0$, and therefore

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

CONSISTENT CLASSIFIER ALGORITHM FOR FINITE \mathcal{F}

From last slide:

$$\Pr \left[\forall f \in \mathcal{F} \cdot \text{err}(f) \leq \text{err}(f, S) + \sqrt{\frac{2 \text{err}(f, S) \ln(|\mathcal{F}|/\delta)}{|S|}} + \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

Since the Consistent Classifier Algorithm returns $\hat{f} \in \mathcal{F}$, we know that

$$\Pr \left[\text{err}(\hat{f}) \leq \text{err}(\hat{f}, S) + \sqrt{\frac{2 \text{err}(\hat{f}, S) \ln(|\mathcal{F}|/\delta)}{|S|}} + \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

By definition of \hat{f} , $\text{err}(\hat{f}, S) = 0$, and therefore

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta.$$

True error of \hat{f} goes to zero as $|S| \rightarrow \infty$ at $O\left(\frac{\log(|\mathcal{F}|/\delta)}{|S|}\right)$ rate. □

Another interpretation:

- ▶ Suppose learning algorithm finds classifier \hat{f} with $\text{err}(\hat{f}, S) = 0$
i.e., a **perfect classifier on the training data!**
(This is possible under realizability assumption.)

GENERALIZATION

Another interpretation:

- ▶ Suppose learning algorithm finds classifier \hat{f} with $\text{err}(\hat{f}, S) = 0$
i.e., a perfect classifier on the training data!
(This is possible under realizability assumption.)
- ▶ How does this perfection generalize to future examples?

GENERALIZATION

Another interpretation:

- ▶ Suppose learning algorithm finds classifier \hat{f} with $\text{err}(\hat{f}, S) = 0$
i.e., a **perfect classifier on the training data!**
(This is possible under realizability assumption.)
- ▶ How does this **perfection** *generalize* to future examples?
- ▶ **Theory says:** with high probability (over random training data),
true error is not much larger than **training error**:

$$\text{err}(\hat{f}) \leq O\left(\frac{\log |\mathcal{F}|}{|S|}\right).$$

Sometimes **true error** is also called **generalization error**.

GENERALIZATION

Another interpretation:

- ▶ Suppose learning algorithm finds classifier \hat{f} with $\text{err}(\hat{f}, S) = 0$
i.e., a **perfect classifier on the training data!**
(This is possible under realizability assumption.)
- ▶ How does this **perfection** *generalize* to future examples?
- ▶ **Theory says:** with high probability (over random training data),
true error is not much larger than **training error**:

$$\text{err}(\hat{f}) \leq O\left(\frac{\log |\mathcal{F}|}{|S|}\right).$$

Sometimes **true error** is also called **generalization error**.

- ▶ Clearly only **reasonable if $\log |\mathcal{F}|$ is finite** and not too large!

INFINITE FUNCTION CLASSES

What about infinite function classes? (e.g., rectangles, linear classifiers)

INFINITE FUNCTION CLASSES

What about infinite function classes? (e.g., rectangles, linear classifiers)

Okay as long as $\#$ **effective behaviors of \mathcal{F} w.r.t. S** is relatively small.

INFINITE FUNCTION CLASSES

What about infinite function classes? (e.g., rectangles, linear classifiers)

Okay as long as **# effective behaviors of \mathcal{F} w.r.t. S** is relatively small.

Let $S := ((\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$. Then

$$\mathcal{F}|_S := \left\{ (f(\mathbf{x}^{(1)}), f(\mathbf{x}^{(2)}), \dots, f(\mathbf{x}^{(n)})) : f \in \mathcal{F} \right\} \subseteq \{0, 1\}^{|S|}$$

(i.e., all the different ways S can be labeled by functions in \mathcal{F}).

INFINITE FUNCTION CLASSES

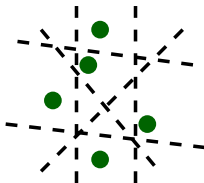
What about infinite function classes? (e.g., rectangles, linear classifiers)

Okay as long as **# effective behaviors of \mathcal{F} w.r.t. S** is relatively small.

Let $S := ((\mathbf{x}^{(1)}, y^{(1)}), (\mathbf{x}^{(2)}, y^{(2)}), \dots, (\mathbf{x}^{(n)}, y^{(n)}))$. Then

$$\mathcal{F}_{|S} := \left\{ (f(\mathbf{x}^{(1)}), f(\mathbf{x}^{(2)}), \dots, f(\mathbf{x}^{(n)})) : f \in \mathcal{F} \right\} \subseteq \{0, 1\}^{|S|}$$

(i.e., all the different ways S can be labeled by functions in \mathcal{F}).



Text

(Some of the ways to label the five points by linear classifiers—there are several more.)

EFFECTIVE BEHAVIORS

What is the size of \mathcal{F}_S ?

EFFECTIVE BEHAVIORS

What is the size of $\mathcal{F}_{|S|}$?

Some possibilities:

- ▶ **Bad situation:** $|\mathcal{F}_{|S|}| = 2^{|S|}$ (all labelings possible).

Function class is too rich for this data set S . Enough functions in \mathcal{F} to perfectly explain all possible labelings (even a random labeling).

EFFECTIVE BEHAVIORS

What is the size of $\mathcal{F}|_S$?

Some possibilities:

- ▶ **Bad situation:** $|\mathcal{F}|_S| = 2^{|S|}$ (all labelings possible).

Function class is too rich for this data set S . Enough functions in \mathcal{F} to perfectly explain all possible labelings (even a random labeling).

- ▶ **Good situation:** $|\mathcal{F}|_S| \leq (c|S|)^v$ for some constant $c > 0$ and non-negative integer v .

Function class \mathcal{F} is limited in capacity to assign labels to points in S .

EFFECTIVE BEHAVIORS

What is the size of $\mathcal{F}_{|S|}$?

Some possibilities:

- ▶ **Bad situation:** $|\mathcal{F}_{|S|}| = 2^{|S|}$ (all labelings possible).

Function class is too rich for this data set S . Enough functions in \mathcal{F} to perfectly explain all possible labelings (even a random labeling).

- ▶ **Good situation:** $|\mathcal{F}_{|S|}| \leq (c|S|)^v$ for some constant $c > 0$ and non-negative integer v .

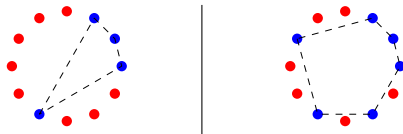
Function class \mathcal{F} is limited in capacity to assign labels to points in S .

We say \mathcal{F} is a **VC class**[†] if, as the number of training data $|S|$ increases, we are eventually in the Good situation (regardless of the actual points in S).

[†]VC = Vapnik-Chervonenkis (1971), same duo who proposed SVMs.

EFFECTIVE BEHAVIORS

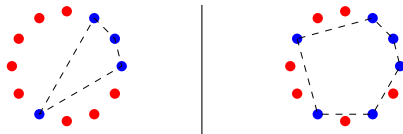
Example: arbitrary convex shapes in \mathbb{R}^2



For any n , there is a set of n points in \mathbb{R}^2 for which convex shapes realize all possible labelings. \rightarrow **Could always be in the Bad situation.**

EFFECTIVE BEHAVIORS

Example: arbitrary convex shapes in \mathbb{R}^2



For any n , there is a set of n points in \mathbb{R}^2 for which convex shapes realize all possible labelings. \rightarrow **Could always be in the Bad situation.**

Example: linear classifiers in \mathbb{R}^2



There are no sets of 4 points in \mathbb{R}^2 where linear classifiers realize all possible labelings. \rightarrow It turns out that $|\mathcal{F}_{|S|}| \leq (c|S|)^3$.

EFFECTIVE BEHAVIORS

Would be great if we could “plug-in” $|\mathcal{F}_S|$ in place of $|\mathcal{F}|$ in guarantee for Consistent Classifier Algorithm:

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2 \ln(|\mathcal{F}|/\delta)}{|S|} \right] \geq 1 - \delta. \quad (\star)$$

EFFECTIVE BEHAVIORS

Would be great if we could “plug-in” $|\mathcal{F}_S|$ in place of $|\mathcal{F}|$ in guarantee for Consistent Classifier Algorithm:

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2 \ln(|\mathcal{F}_S|/\delta)}{|S|} \right] \geq 1 - \delta. \quad (\star)$$

EFFECTIVE BEHAVIORS

Would be great if we could “plug-in” $|\mathcal{F}_{|S|}|$ in place of $|\mathcal{F}|$ in guarantee for Consistent Classifier Algorithm:

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2 \ln(|\mathcal{F}_{|S|}|/\delta)}{|S|} \right] \geq 1 - \delta. \quad (\star)$$

► **Bad situation:** $|\mathcal{F}_{|S|}| = 2^{|S|}$ (all labelings possible). Useless in (\star) :

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2|S| \ln(2) + 2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

EFFECTIVE BEHAVIORS

Would be great if we could “plug-in” $|\mathcal{F}_{|S|}|$ in place of $|\mathcal{F}|$ in guarantee for Consistent Classifier Algorithm:

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2 \ln(|\mathcal{F}_{|S|}|/\delta)}{|S|} \right] \geq 1 - \delta. \quad (\star)$$

- **Bad situation:** $|\mathcal{F}_{|S|}| = 2^{|S|}$ (all labelings possible). Useless in (\star) :

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2|S| \ln(2) + 2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

- **Good situation:** $|\mathcal{F}_{|S|}| \leq (c|S|)^v$. Bound in (\star) becomes

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2v \ln(c|S|) + 2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

EFFECTIVE BEHAVIORS

Would be great if we could “plug-in” $|\mathcal{F}_{|S|}|$ in place of $|\mathcal{F}|$ in guarantee for Consistent Classifier Algorithm:

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2 \ln(|\mathcal{F}_{|S|}|/\delta)}{|S|} \right] \geq 1 - \delta. \quad (\star)$$

- **Bad situation:** $|\mathcal{F}_{|S|}| = 2^{|S|}$ (all labelings possible). Useless in (\star) :

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2|S| \ln(2) + 2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

- **Good situation:** $|\mathcal{F}_{|S|}| \leq (c|S|)^v$. Bound in (\star) becomes

$$\Pr \left[\text{err}(\hat{f}) \leq \frac{2v \ln(c|S|) + 2 \ln(1/\delta)}{|S|} \right] \geq 1 - \delta.$$

This straight-up “plugging-in” isn't technically legal,
but different argument implies something like (\star) is true!

RECAP AND FINAL REMARKS

- ▶ The Consistent Classifier Algorithm returns $\hat{f} \in \mathcal{F}$ with $\text{err}(\hat{f}) \rightarrow 0$ as $|S| \rightarrow \infty$ with high probability, provided that:
 - ▶ labels are realized by some $f^* \in \mathcal{F}$;
 - ▶ $|\mathcal{F}|$ is finite, or $|\mathcal{F}_{|S}|$ only grows polynomially with $|S|$.

RECAP AND FINAL REMARKS

- ▶ The Consistent Classifier Algorithm returns $\hat{f} \in \mathcal{F}$ with $\text{err}(\hat{f}) \rightarrow 0$ as $|S| \rightarrow \infty$ with high probability, provided that:
 - ▶ labels are realized by some $f^* \in \mathcal{F}$;
 - ▶ $|\mathcal{F}|$ is finite, or $|\mathcal{F}_{|S}|$ only grows polynomially with $|S|$.
- ▶ Guarantees depends on **complexity of function class \mathcal{F}** (either cardinality or effective number of behaviors).

RECAP AND FINAL REMARKS

- ▶ The Consistent Classifier Algorithm returns $\hat{f} \in \mathcal{F}$ with $\text{err}(\hat{f}) \rightarrow 0$ as $|S| \rightarrow \infty$ with high probability, provided that:
 - ▶ labels are realized by some $f^* \in \mathcal{F}$;
 - ▶ $|\mathcal{F}|$ is finite, or $|\mathcal{F}_S|$ only grows polynomially with $|S|$.
- ▶ Guarantees depends on **complexity of function class \mathcal{F}** (either cardinality or effective number of behaviors).
- ▶ *Without realizability assumption*, essentially same argument applies to give a different guarantee.

Using ERM: $\hat{f} := \arg \min_{f \in \mathcal{F}} \text{err}(f, S)$, with high probability, *excess error*

**get close to
best classifier**

$$\text{err}(\hat{f}) - \min_{f \in \mathcal{F}} \text{err}(f)$$

goes to zero as $|S|$ increases under same complexity conditions.

CROSS VALIDATION

MODEL SELECTION

Objective

- ▶ Often necessary to consider many different models for a given problem (e.g., class conditional distributions in generative model classifiers, features in linear classifiers, kernel in kernelized classifiers).
- ▶ Sometimes “model” simply means particular setting of **hyper-parameters** (e.g., k in k -NN, λ in soft-margin SVM, number of nodes in decision tree).

Terminology

The problem of choosing a good model is called **model selection**.

EXAMPLE: SVM WITH GAUSSIAN KERNEL

Soft-margin SVM with Gaussian kernel

- ▶ Models indexed by regularization parameter λ and Gaussian kernel bandwidth $h > 0$:

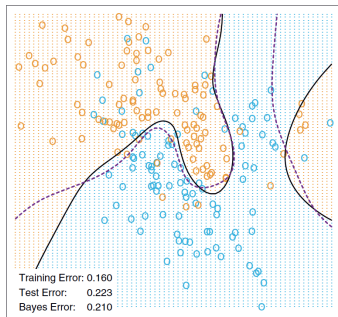
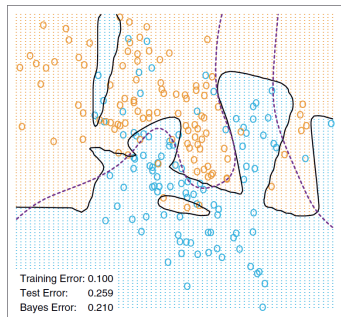
$$K(\mathbf{x}, \tilde{\mathbf{x}}) = \exp\left(-\frac{\|\mathbf{x} - \tilde{\mathbf{x}}\|_2^2}{2h}\right).$$

- ▶ Goal is to find setting of (λ, h) for which we can expect small true (generalization) error.

Naïve approach

- ▶ Also minimize over (λ, h) in SVM optimization problem.
- ▶ Leads to **overfitting**: resulting SVM classifier adapts too closely to specific properties of the training data, rather than underlying distribution.

OVERFITTING: ILLUSTRATION



- ▶ Classifier in this example has bandwidth parameter σ (similar to Gaussian kernel bandwidth).
- ▶ Small $\sigma \rightarrow$ permits curve with sharp bends
- ▶ Large $\sigma \rightarrow$ smoother curve.

MODEL SELECTION BY HOLD-OUT VALIDATION

(Henceforth, use \mathbf{h} to denote particular setting of hyper-parameters / model choice.)

Hold-out validation

Model selection:

1. Randomly split data into three sets: **training**, **validation**, and **test** data.

Training	Validation	Test
----------	------------	------

2. Train classifier $\hat{f}_{\mathbf{h}}$ on **Training** data for different values of \mathbf{h} .
3. Compute **Validation** ("hold-out") error for each $\hat{f}_{\mathbf{h}}$: $\text{err}(\hat{f}_{\mathbf{h}}, \text{Validation})$.
4. Selection: $\hat{\mathbf{h}}$ = value of \mathbf{h} with lowest **Validation** error.
5. Train classifier \hat{f} using $\hat{\mathbf{h}}$ with **Training** + **Validation** data.

Model assessment:

6. Finally: estimate the error of \hat{f} using **test** data.

MAIN IDEA BEHIND HOLD-OUT VALIDATION

Training	Validation	Test
----------	------------	------

Classifier \hat{f}_h trained on Training data $\rightarrow \text{err}(\hat{f}_h, \text{Validation})$.

Training + Validation	Test
-----------------------	------

Classifier \hat{f}_h trained on Training + Validation data $\rightarrow \text{err}(\hat{f}_h, \text{Test})$.

MAIN IDEA BEHIND HOLD-OUT VALIDATION

Training	Validation	Test
----------	------------	------

Classifier \hat{f}_h trained on Training data $\rightarrow \text{err}(\hat{f}_h, \text{Validation})$.

Training + Validation	Test
-----------------------	------

Classifier \hat{f}_h trained on Training + Validation data $\rightarrow \text{err}(\hat{f}_h, \text{Test})$.

The hope is that these quantities are similar!

MAIN IDEA BEHIND HOLD-OUT VALIDATION

Training	Validation	Test
----------	------------	------

Classifier \hat{f}_h trained on Training data $\rightarrow \text{err}(\hat{f}_h, \text{Validation})$.

Training + Validation	Test
-----------------------	------

Classifier \hat{f}_h trained on Training + Validation data $\rightarrow \text{err}(\hat{f}_h, \text{Test})$.

for a lot ~
it holds

The hope is that these quantities are similar!

(Making this rigorous is actually rather tricky.)

BEYOND SIMPLE HOLD-OUT VALIDATION

Standard hold-out:

Training	Validation	Test
----------	------------	------

Classifier \hat{f}_h trained on Training data $\rightarrow \text{err}(\hat{f}_h, \text{Validation})$.

BEYOND SIMPLE HOLD-OUT VALIDATION

Standard hold-out:

Training	Validation	Test
----------	------------	------

Classifier \hat{f}_h trained on Training data $\rightarrow \text{err}(\hat{f}_h, \text{Validation})$.

Could also:

- ▶ train \hat{f}_h using Validation data, and
- ▶ evaluate \hat{f}_h using Training data.

Training	Validation	Test
----------	------------	------

Classifier \hat{f}_h trained on Validation data $\rightarrow \text{err}(\hat{f}_h, \text{Training})$.

BEYOND SIMPLE HOLD-OUT VALIDATION

Standard hold-out:

Training	Validation	Test
----------	------------	------

Classifier \hat{f}_h trained on Training data $\rightarrow \text{err}(\hat{f}_h, \text{Validation})$.

Could also:

- ▶ train \hat{f}_h using Validation data, and
- ▶ evaluate \hat{f}_h using Training data.

Training	Validation	Test
----------	------------	------

Classifier \hat{f}_h trained on Validation data $\rightarrow \text{err}(\hat{f}_h, \text{Training})$.

Idea: Do both, and average results as overall validation error for h .

MODEL SELECTION BY K -FOLD CROSS VALIDATION

Model selection:

1. Set aside some **test** data.
2. Of remaining data, split into K parts ("folds") S_1, S_2, \dots, S_K .
3. For each value of h :
 - ▶ For each $k \in \{1, 2, \dots, K\}$:

- ▶ Train classifier $\hat{f}_{h,k}$ using all S_i except S_k .
- ▶ Evaluate classifier $\hat{f}_{h,k}$ using S_k : $\text{err}(\hat{f}_{h,k}, S_k)$

Example: $K = 5$ and $k = 4$

Training	Training	Training	Validation	Training
----------	----------	----------	------------	----------

- ▶ Cross-validation error for h : $\frac{1}{K} \sum_{k=1}^K \text{err}(\hat{f}_{h,k}, S_k)$.

4. Select the value \hat{h} with lowest cross-validation error.

5. Train classifier \hat{f} using selected \hat{h} with all S_1, S_2, \dots, S_K .

Model assessment:

6. Finally: estimate the error of \hat{f} using **test** data.

note after the classifier was selected..

we need to train on all data to get classifier

HOW TO CHOOSE K ?

Argument for small K

Better simulates “variation” between different training samples drawn from underlying distribution.

$$K = 2$$

Training	Validation
Validation	Training

$$K = 4$$

Validation	Training	Training	Training
Training	Validation	Training	Training
Training	Training	Validation	Training
Training	Training	Training	Validation

HOW TO CHOOSE K ?

Argument for small K

Better simulates “variation” between different training samples drawn from underlying distribution.

$K = 2$

Training	Validation
Validation	Training

$K = 4$

Validation	Training	Training	Training
Training	Validation	Training	Training
Training	Training	Validation	Training
Training	Training	Training	Validation

Argument for large K

Some learning algorithms exhibit *phase transition* behavior (e.g., output is complete rubbish until sample size sufficiently large).
Using large K best simulates training on **all** data (except **test**, of course).

HOW TO CHOOSE K ?

Argument for small K

Better simulates “variation” between different training samples drawn from underlying distribution.

$K = 2$

Training	Validation
Validation	Training

$K = 4$

Validation	Training	Training	Training
Training	Validation	Training	Training
Training	Training	Validation	Training
Training	Training	Training	Validation

Argument for large K

Some learning algorithms exhibit *phase transition* behavior

(e.g., output is complete rubbish until sample size sufficiently large).

Using large K best simulates training on **all** data (except **test**, of course).

In practice: usually $K = 5$ or $K = 10$.

- ▶ **Model selection:** goal is to pick best model (e.g., features, kernels, hyper-parameter settings) to achieve low true error.

- ▶ **Model selection:** goal is to pick best model (e.g., features, kernels, hyper-parameter settings) to achieve low true error.
- ▶ **Two common methods:** hold-out validation and K -fold cross validation (with $K = 5$ or $K = 10$).

- ▶ **Model selection:** goal is to pick best model (e.g., features, kernels, hyper-parameter settings) to achieve low true error.
- ▶ **Two common methods:** hold-out validation and K -fold cross validation (with $K = 5$ or $K = 10$).
- ▶ **Caution:** considering *too many* different models can lead to overfitting, even with hold-out / cross-validation.

- ▶ **Model selection:** goal is to pick best model (e.g., features, kernels, hyper-parameter settings) to achieve low true error.
- ▶ **Two common methods:** hold-out validation and K -fold cross validation (with $K = 5$ or $K = 10$).
- ▶ **Caution:** considering *too many* different models can lead to overfitting, even with hold-out / cross-validation.

(Sometimes “*averaging*” the models in some way can help.)