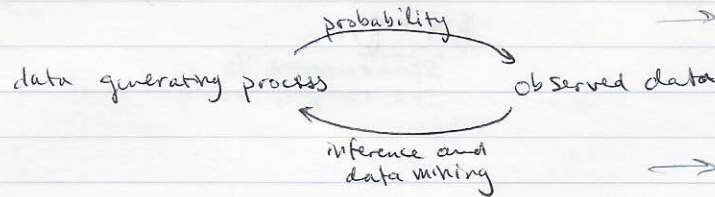


# DME

## final exam notes

### intro lecture

- there's, like, a lot of data
- data mining  $\approx$  data analysis  $\approx$  data science
- data analysis as statistical inference



given a data generating process, what are the properties of the outcomes?

given the outcomes (data), what can we say about the process that generated them?

### data analysis process

- get (raw) data (lecture 5)
  - understand sampling processes
  - any biases?
  - feedback loops between data analysis and collection?
- exploratory data analysis (lectures 1-3)
  - become familiar with data
  - spot unexpected properties
  - anomalies, outliers, missing data?
- prep data for further analysis
  - merge data sets, reformat
  - select/exclude data
  - provide clear rationale for selection/execution
- build and fit model (lecture 4)
  - generalization is the goal
  - choice of evaluation metric
  - choice of hyperparameters
- summarize, visualize results
- deploy the product / communicate findings

start again

presentations, mini-project

### first steps in exploratory data analysis

- distributions of single variables

- numerical summaries

$$\vec{m} = \sum_{i=1}^n \frac{1}{n} \mathbf{1}_n = \left[ \frac{1}{n} \sum_{i=1}^n x_i, \dots, \frac{1}{n} \sum_{i=1}^n 1 \right]$$

- location
- mean:  $m = \frac{1}{n} \sum_{i=1}^n x_i$ ,  $E(x) = \int x p(x) dx$
  - median: robust; if some outlier exists, only changes median to value of neighboring point
  - trimmed mean:  $\frac{1}{n-2k} \sum_{i=k+1}^{n-k} x_{(i)} \Rightarrow k=0$  normal mean,  $k \rightarrow \frac{n}{2}$  median

- scale
- variance:  $v = \frac{1}{n} \sum_{i=1}^n (x_i - m)^2 = E(x^2) - E(x)^2$
  - median absolute deviation is more robust:  $MAD = \text{median}(|x_i - \text{median}(x_i)|)$
  - IQR (also robust): range of middle 50% of data:  $IQR = x_{(13/4n)} - x_{(1/4n)}$

shape (skew)

- skew  $(x) = E \left[ \left( \frac{x - \mu}{\sigma} \right)^3 \right] \Rightarrow = 0$  if symmetric around mean  
 $> 0$  if long tail to the right
- Galton's measure of skewness:  $\frac{(Q_3 - Q_2) - (Q_2 - Q_1)}{Q_3 - Q_1}$



- kurt(x) =  $E \left[ \left( \frac{x-\mu}{\sigma} \right)^4 \right] \Rightarrow = 3$  for Gaussian

- excess kurtosis (x) = kurt(x) - 3  
spread around  $\frac{3}{4}$       spread around  $\frac{1}{4}$

- robust kurtosis (x) =  $\frac{(q_{7/8} - q_{5/8}) + (q_{3/8} - q_{1/8})}{q_{3/4} - q_{1/4}}$   
spread around  $\frac{1}{2}$   
 (x2 range of two individual above)

- graphs

• histogram

- binning:  $B_1 = [L, L+h)$ ,  $B_2 = [L+h, L+2h)$ , ...,  $B_k = [L+(k-1)h, L+kh)$

• kernel density estimate

- fixes problems w/ histogram: doesn't depend on starting point, smooth

- boxcar: counts points  $\pm h/2$ :

$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} \Pi_h(x - x_i)$  where  $\Pi_h(x) = \begin{cases} 1 & \text{if } x \in [-\frac{h}{2}, \frac{h}{2}] \\ 0 & \text{otherwise} \end{cases}$

- Gaussian kernel:

$\hat{p}(x) = \frac{1}{n} \sum_{i=1}^n K_h(x - x_i)$  where  $K_h(x) = \frac{1}{\sqrt{2\pi}h} \exp\left(-\frac{x^2}{2h^2}\right)$

- boxplot: median,  $Q_1/Q_3$ , and  $Q_1 - 1.5 IQR$ ,  $Q_3 + 1.5 IQR$

• joint distributions of two variables

→ numerical summaries

also,  $C = E[(x-\mu)(x-\mu)^T]$

•  $\text{cov}(x, y) = E[(x-\mu_x)(y-\mu_y)] = E[xy] - \mu_x \mu_y$

$\text{cov}(ax+b, y) = a \text{cov}(x, y)$

• correlation  $\rho(x, y) = \frac{\text{cov}(x, y)}{\sqrt{V(x)V(y)}} = E\left[\left(\frac{x-\mu_x}{\sigma_x}\right)\left(\frac{y-\mu_y}{\sigma_y}\right)\right]$

• multiplying by constant:

$V(Ax+b) = AC A^T$  where  $C = V(x)$

• correlation matrix K: use the linear transformation  $D^{-1/2}(x-\mu)$  where D contains the diagonal elements of C, i.e.,  $D = \text{diag}(V(x_1), \dots, V(x_d))$

$K = D^{-1/2} C D^{-1/2}$

- by construction, all diagonal elements are 1

• nonlinear relationships:

- correlation:  $\rho(g(x), g(y)) = \frac{\text{cov}(g(x), g(y))}{\sqrt{V(g(x))V(g(y))}}$

• simple preprocessing

- simple outlier detection: see if outside range  $[Q_1 - k IQR, Q_3 + k IQR]$  where  $k=1.5$  usually

- data standardization

• centering:  $X = \tilde{X} H_n$  where X is centered data,  $\tilde{X}$  is uncentered data,  $H_n = I_n - \frac{1}{n} \mathbf{1}_n \mathbf{1}_n^T$

or:  $X = (\tilde{x}_1 - \underline{m}), \dots, (\tilde{x}_n - \underline{m})$   
 $= \tilde{X} - (\underline{m}, \dots, \underline{m})$

$\tilde{X} H_n$  sums values of a  
 $= \tilde{X} - \frac{1}{n} \tilde{X} \mathbf{1}_n \mathbf{1}_n^T$   
 $= \tilde{X} - \tilde{X} \frac{1}{n} \mathbf{1}_n^T \mathbf{1}_n$

$H_n = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{bmatrix} - \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}$

and  $\underline{m} = \tilde{X} \frac{1}{n} \mathbf{1}_n$

- scaling to unit variance:  $\hat{C} = \frac{1}{n} \tilde{X} H_n \tilde{X}^T$



# DME

## final exam notes

### principal component analysis,

- assume for this lecture that data has been centered
- PCA by sequential variance maximization
  - first PCA direction
    - unit vector  $\underline{w}_1$  for which the projected data  $\underline{w}_1^T \underline{x}_i$  is maximally variable

$$\max_{\underline{w}_1} \underline{w}_1^T C \underline{w}_1$$

$$\text{s.t. } \|\underline{w}_1\| = 1$$

do this because want to max variance  $V(\underline{w}_1^T \underline{x})$

note for later:  $C = \frac{1}{n} \underline{X} \underline{X}^T$

- can be solved in closed form from  $C = \underline{U} \underline{\Lambda} \underline{U}^T$

where  $\underline{U}$  is an orthogonal matrix,  $\underline{\Lambda}$  is a diagonal matrix with eigenvalues  $\lambda_i \geq 0$   
and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$

$$\text{let } \underline{w}_1 = \underline{U} \underline{a}$$

$$\text{then } \underline{w}_1^T C \underline{w}_1 = \underline{a}^T \underline{U}^T \underline{U} \underline{\Lambda} \underline{U}^T \underline{U} \underline{a} = \underline{a}^T \underline{\Lambda} \underline{a} = \sum_{i=1}^d a_i^2 \lambda_i$$

$$\text{and constraint becomes } \|\underline{w}_1\|^2 = \underline{w}_1^T \underline{w}_1$$

$$= \underline{a}^T \underline{U}^T \underline{U} \underline{a} = \underline{a}^T \underline{a} = \sum_{i=1}^d a_i^2 = 1$$

so formulation is now:

$$\max_{a_1, \dots, a_d} \sum_{i=1}^d a_i^2 \lambda_i$$

$$\text{s.t. } \sum_{i=1}^d a_i^2 = 1$$

to solve, just set  $a_1 = 1$  and the remaining  $a_i$  to zero

$$\underline{w}_1 = \underline{U} \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \underline{u}_1$$

score

- first principal component (not PC direction) is:  $z_1 = \underline{w}_1^T \underline{x}$

- variance of  $z_1$  is  $\lambda_1$

$$= \underline{u}_1^T \underline{x} \rightarrow \text{also by } s_1^T \underline{v}_1 \text{, etc.}$$

- all PC scores given by  $\underline{w}_i^T \underline{x} = \underline{u}_i^T \underline{x}$

- subsequent PC directions:

$$\max_{\underline{w}_2} \underline{w}_2^T C \underline{w}_2$$

$$\text{s.t. } \|\underline{w}_2\| = 1$$

$$\Rightarrow \max_{b_1, \dots, b_d} \sum_{i=1}^d b_i^2 \lambda_i$$

$$\text{s.t. } \sum_{i=1}^d b_i^2 = 1$$

letting  $\underline{w}_2 = \underline{U} \underline{b}$

$\underline{w}_2^T \underline{w}_1 = 0 \rightarrow \underline{w}_2^T \underline{x}$  needs to reveal something "new" about the data

$$\underline{b}_1 = 0$$

$$\text{since } \underline{w}_2^T \underline{w}_1 = \underline{b}^T \underline{U}^T \underline{u}_1$$

$$= \underline{b}^T \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$$= b_1$$

$$= 0 \text{ since}$$

$$\underline{w}_2^T \underline{w}_1 = 0$$

$$\Rightarrow \max_{b_2, \dots, b_d} \sum_{i=2}^d b_i^2 \lambda_i$$

$$\text{s.t. } \sum_{i=2}^d b_i^2 = 1$$

this is structurally similar to before; we get

$$\underline{w}_2 = \underline{U} \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} = \underline{u}_2$$

- Subsequent PC directions (cont'd)

• so in general:

$$\max_{\underline{w}_m} \underline{w}_m^T \underline{C} \underline{w}_m$$

$$\text{s.t. } \|\underline{w}_m\| = 1$$

$$\underline{w}_m^T \underline{w}_i = 0 \text{ for } i=1, \dots, m-1$$

• total variance explained of the  $\underline{z}_m, m=1, \dots, k$  is

$$\sum_{m=1}^k V(\underline{z}_m) = \sum_{m=1}^k \lambda_m$$

• PCA by simultaneous variance maximization

- Solves the problem:

$$\max_{\underline{w}_1, \dots, \underline{w}_k} \sum_{i=1}^k \underline{w}_i^T \underline{C} \underline{w}_i$$

$$\text{s.t. } \|\underline{w}_i\| = 1 \quad i=1, \dots, k$$

$$\underline{w}_i^T \underline{w}_j = 0 \quad i \neq j$$

It's interesting that the sequential (greedy) approach yields the same answer since greedy approaches usually don't

• PCA by maximization of approximation error

- have a set of  $k$  orthonormal vectors  $\underline{w}_1, \dots, \underline{w}_k$

$$\underline{P} = \sum_{i=1}^k \underline{w}_i \underline{w}_i^T = \underline{W}_k \underline{W}_k^T \quad \text{where } \underline{W}_k = (\underline{w}_1, \dots, \underline{w}_k)$$

where  $\underline{P}$  projects any vector onto its subspace

$$\hat{\underline{x}}_i = \underline{P} \underline{x}_i \quad , \quad \hat{\underline{x}} = \underline{P} \underline{x} = \sum_{i=1}^k \underline{w}_i \underline{w}_i^T \underline{x}$$

ask: which subspace yields the smallest expected approximation error?

$$\min_{\underline{w}_1, \dots, \underline{w}_k} \mathbb{E} \left\| \underline{x} - \sum_{i=1}^k \underline{w}_i \underline{w}_i^T \underline{x} \right\|^2$$

$$\text{s.t. } \|\underline{w}_i\| = 1 \quad i=1, \dots, k$$

$$\underline{w}_i^T \underline{w}_j = 0 \quad i \neq j$$

after some algebra:

$$\mathbb{E} \left\| \underline{x} - \underline{W}_k \underline{W}_k^T \underline{x} \right\|^2 = \mathbb{E}(\underline{x}^T \underline{x}) - \sum_{i=1}^k \underline{w}_i^T \underline{C} \underline{w}_i \quad \rightarrow \text{so minimizing approx error is equivalent to max } \sum_{i=1}^k \underline{w}_i^T \underline{C} \underline{w}_i$$

(one minus) constant

- computing fraction of variance explained:

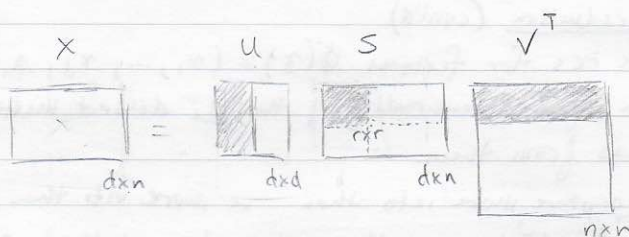
$$\frac{\mathbb{E} \left\| \underline{x} - \underline{U}_k \underline{U}_k^T \underline{x} \right\|^2}{\mathbb{E}(\underline{x}^T \underline{x})} = 1 - \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} = 1 - [\text{fraction of variance explained}]$$



## Principal component analysis (contd)

• SVD:

$$X = U S V^T$$



So now it's an approximation

$$X = \sum_{i=1}^r s_i \underline{u}_i \underline{v}_i^T$$

$$= \sum \underbrace{u_i v_i s_i}_{\text{PC directions (same } u_i \text{ as before)}} \quad \text{PC scores } (z_i)$$

$$= \sum_{i=1}^k u_i z_i^T$$

- if the Gram matrix is defined as  $G = X^T X$ , we can find its eigenvalue decomposition as

$$G = (USV^T)^T (USV) = \dots = V \Sigma V^T \text{ where } \Sigma = SS^T \rightarrow \text{eigenvalues } s_i^2$$

So we have  $s_i$ 's

and  $V$ , which is what we need to compute the PC scores

we can therefore compute the PC scores of  $X$  from its Gram matrix of centered data, without first computing the PC directions

$$Z = \sqrt{\Sigma_k} V_k^T$$

$$= S_k V_k^T$$

also,  $Z = U_K^T X$   
but we don't have  $U_K$   
if we only have the Gram  
matrix

dimensionality reduction

• can compute  $G$  from uncentered  $\tilde{G}$ :

$$G = X^T X = H_n^T \tilde{X}^T \tilde{X} H_n = H_n^T \tilde{G} H_n$$

(need to use fact that  $H_n$  is symmetric)

So just need to "double center" first

- can compute  $G$  from  $\Delta$ , where  $\Delta$  has elements  $\delta_{ij}^2 = \|\bar{x}_i - \bar{x}_j\|^2 = \|\bar{x}_i\|^2 + \|\bar{x}_j\|^2 - 2\bar{x}_i^T \bar{x}_j$

Since  $\|x_i\|^2$  is a constant along row  $i$   
and  $\|x_j\|^2$  is a constant along column  $j$ ,

we can double-center and multiply by  $-\frac{1}{2}$  (to recover  $\mathbf{x}_i^\top \mathbf{x}_j$ ):

$$G = -\frac{1}{2} H \Delta H$$



## dimensionality reduction (cont'd)

idea: compute PCs for features  $\phi(\underline{x}) = (x_1, \dots, x_d, x_1 x_2, \dots, x_1 x_d, \dots, x_d x_d)^T$

kernel  
PCA

- the much higher dimensionality of the  $\phi_i$  doesn't matter as long as we only compute  $k$  principal components from them.
- PCs computed from  $\phi_i$  must capture more info than PCs captured from  $\underline{x}_i$
- kernel trick: don't actually need to know individual  $\phi_i$ , only the inner products  $\phi_i^T \phi_j = \phi(\underline{x}_i)^T \phi(\underline{x}_j)$   
 $\Rightarrow$  then we can just compute the PC scores from the Gram matrix of  $\Phi$ :  
 $(\tilde{G})_{ij} = \phi_i^T \phi_j = \phi(\underline{x}_i)^T \phi(\underline{x}_j) = k(\underline{x}_i, \underline{x}_j)$

- example kernels:

- polynomial:  $k(\underline{x}, \underline{x}') = (\underline{x}^T \underline{x}')^a$
- Gaussian:  $k(\underline{x}, \underline{x}') = \exp\left(-\frac{\|\underline{x} - \underline{x}'\|^2}{2\sigma^2}\right)$

- to compute PCs, proceed exactly as before:

$$G = H_n \tilde{G} H_n \Rightarrow Z = \sqrt{\tilde{\Sigma}_k} V_k^T$$

• multidimensional scaling (MDS)

- umbrella term for several methods that operate on dissimilarities  $\delta_{ij}$  (one simple example is Euclidean distance)
- metric MDS: numerical values of the dissimilarities are assumed to carry info  
 (contrast this w/ nonmetric MDS, where only the rank-order of the dissimilarities matter)

$\rightarrow$  "stress": essentially a cost function

$$\min_{z_1, \dots, z_n} \sum_{i < j} w_{ij} (\|\underline{z}_i - \underline{z}_j\| - \delta_{ij})^2 \quad \text{i.e., distance between points we create a distance according to distance metric}$$

- goal is to find  $n$  points  $\underline{z} \in \mathbb{R}^k$  that solve this,
- $\|\underline{z}_i - \underline{z}_j\|$  is the Euclidean distance between  $\underline{z}_i$  and  $\underline{z}_j$
- $k$  usually set to 2 for visualization
- usually solved by gradient descent
- weights  $w_{ij} \geq 0$  are specified by the user
- nonmetric MDS: only relation between the  $\delta_{ij}$  are assumed to matter  
 optimization modified to:

$$\min_{z_1, \dots, z_n, f} \sum_{i < j} w_{ij} (\|\underline{z}_i - \underline{z}_j\| - f(\delta_{ij}))^2 \quad \text{where } f \text{ is a monotonic function that converts the dissimilarities into distances}$$

- typically solved by iterating between optimization w.r.t. the  $\underline{z}_i$  and w.r.t.  $f$ , which can be done by regression
- classical MDS: assumes dissimilarities  $\delta_{ij}$  are (squared) Euclidean distances
- like before,
  1. compute hypothetical Gram matrix  $G'$ :  $G' = -\frac{1}{2} H_n \Delta H_n$
  2. compute top  $k$  eigenvectors  $\underline{v}_k \in \mathbb{R}^n$  of  $G$  and form matrices  $\tilde{\Sigma}_k$  and  $V_k$
  3. compute PC scores as  $Z = \sqrt{\tilde{\Sigma}_k} V_k^T$
- nice in that it produces nested solutions
- problem:  $\Delta$  is symmetric but not necessarily PSD
- solution: choose  $k$  small enough that all eigenvalues contained in  $\tilde{\Sigma}_k$  are positive



## dimensionality reduction (cont'd)

### • isomap

- uses geodesic distance, which is when dissimilarity is measured by the shortest distance between two points only when allowed to travel on the data manifold from one neighboring datapoint to the next
- with disconnected parts of the graph, the isomap is often applied to each component separately

Neighborhood of a data-point can be taken to be  $m$ -nearest neighbors or all points w/ in a Euclidean distance

## performance evaluation in predictive modeling

### • prediction and training loss

$$J(h) = \mathbb{E}_{x,y} [L(\hat{y}, y)] = \mathbb{E}_{x,y} [L(h(x), y)]$$

- optimize:  $\min_h J(h)$

### - three difficulties:

1. can't compute expectation over  $(x, y)$  analytically
2.  $L$  may be hard to evaluate
3. minimizing  $L$  w.r.t. a function is generally difficult

### - training loss

- approx training loss w/ sample average

$$J(h) \approx \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i)$$

- Samples called  $\mathcal{D}^{\text{train}}$

- indicate model family and hyperparams:  $h(x) \equiv h_{\underline{x}}(x; \theta)$  where  $\underline{x}$  is a vector of hyperparams indicating model family and some tuning params associated with it;  $\theta$  is usually weights

- work with training loss function

$$J_{\underline{x}}(\theta) = \frac{1}{n} \sum_{i=1}^n L(h_{\underline{x}}(x_i; \theta), y_i)$$

### • generalization performance

- algorithm that turns training data  $\mathcal{D}^{\text{train}}$  into a prediction function  $\hat{h}$  by  $A$ :

$$\hat{h} = A(\mathcal{D}^{\text{train}})$$

so:

$$\tilde{J}(A) = \mathbb{E}_{\mathcal{D}^{\text{train}}} [J(\hat{h})] = \mathbb{E}_{\mathcal{D}^{\text{train}}} [J(A(\mathcal{D}^{\text{train}}))]$$

$A$  is more general than  $\hat{h}$

- overfitting: a model has been overfitted to the training data if reducing its complexity reduces the (expected) prediction loss.
- underfitting: <sup>when</sup> increasing flexibility  $\Rightarrow$  decreased (expected) prediction loss
- variability of prediction loss increases with the flexibility of the model
- bias is the difference between training loss and prediction loss as  $n \rightarrow \infty$

### • estimating generalization performance

- hold-out approach:  $\hat{h} = A(\mathcal{D}^{\text{train}})$

validation error:  $\hat{J}(\hat{h}; \tilde{\mathcal{D}}) = \frac{1}{\tilde{n}} \sum_{i=1}^{\tilde{n}} L(\hat{h}(\tilde{x}_i), \tilde{y}_i)$  where  $\tilde{\mathcal{D}}$  is validation set (and tildes denote validation-related variables)

useful to stratify



## performance evaluation in predictive modeling (cont'd)

### • estimating generalization performance (cont'd)

#### - cross validation:

- split data into  $K$  subsets ("folds")  $D_1, \dots, D_K$ , perhaps w/ stratification

$$\text{do for } k=1, \dots, K \left\{ \begin{array}{l} D_k^{\text{train}} = \bigcup_{i \neq k} D_i, \quad D_k^{\text{val}} = D_k \\ \hat{h}_k = A(D_k^{\text{train}}) \\ \hat{J}_k = \hat{J}(\hat{h}_k; D_k^{\text{val}}) \end{array} \right.$$

$$CV = \frac{1}{K} \sum_{k=1}^K \hat{J}_k \Rightarrow \text{estimate of prediction performance } \bar{J}(A) \text{ as opposed to } \hat{J} \text{ from before}$$

- can estimate variance of cv-score:  $V(CV) \approx \frac{1}{K} V(\hat{J})$

$$V(\hat{J}) \approx \frac{1}{K} \sum_{k=1}^K (\hat{J}_k - CV)^2$$

### - hyperparameter selection and performance evaluation

#### • two main approaches:

1. hold-out to select hyperparams, hold-out for performance eval
2. CV " " " " (again)

#### • Approach for "two-times hold-out"

1. Split off test data  $D^{\text{test}}$
2. Of remaining, split in  $D^{\text{train}}$  and  $D^{\text{val}}$
3. for all  $\lambda$ :  
 $\hat{h}_\lambda = A_\lambda(D^{\text{train}})$
4. choose best  $\lambda$  as  $\hat{\lambda} = \arg\min PL(\lambda)$  where  $PL(\lambda) = \hat{J}(\hat{h}_\lambda; D^{\text{val}})$
5. re-estimate model w/ best  $\hat{\lambda}$  on all non-test data:  
 $\hat{h} = A_{\hat{\lambda}}(D^{\text{train}} \cup D^{\text{val}})$

note: CV-score is typically an optimistic estimate of prediction loss b/c hyperparams are chosen to minimize it

6. evaluate on test data:

$$\hat{J} = \hat{J}(\hat{h}; D^{\text{test}})$$

7. [optional] re-estimate  $\hat{h}$  using all data for final prediction function  $\hat{h}$

$$\hat{h}(x) = A_{\hat{\lambda}}(D)$$

- approach for cross-validation and hold-out: similar to above