

EE445 Mod2-Lec2: Least Squares Data Fitting

References:

- [VMLS]: Chapter 13

Outline

- least squares data fitting including feature transformation
- validation and generalization
- feature engineering

Least Squares Data Fitting: What is it?

- In [Mod2-L1] we saw the "least squares approximate solution" to $Ax = b$
- Data fitting is one of the most important applications of least squares
- **Goal:** find a mathematical model, or an approximate model, of some relation, given some observed data say from experiments or real-world phenomena
- **Set-up:** Consider $x \in \mathbb{R}^n$ and $y \in \mathbb{R}$. Suppose we **model** x and y as being approximately related by $f : \mathbb{R}^n \rightarrow \mathbb{R}$:

$$y \approx f(x)$$

- **Components:**
 - ▶ feature vector (independent variables): x
 - ▶ outcome (response/dependent variable): y
- **Challenge:** We do not know f and need to approximate it

Data & Model

- Consider m feature vectors and corresponding scalars:

$$x^{(1)}, \dots, x^{(m)} \in \mathbb{R}^n, \quad \text{and} \quad y^{(1)}, \dots, y^{(m)} \in \mathbb{R}$$

- **Data pair** (observations or samples):

- **Model:**

- **prediction:**

Class of Linear Models

- **Linear Models:** consider p **basis functions** or **feature mappings** $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$ and *model parameters* θ_i :
 - the feature mappings f_i are chosen based on our prior knowledge of the problem
 - Once they are chosen, they are a known part of the problem of data fitting.
 - The remaining unknown is the model parameters $\theta = (\theta_1, \dots, \theta_p)$
 - **Goal:**

Prediction Error

- Consider a fixed value of θ and let $f(x) = \sum_{j=1}^p \theta_j f_j(x)$
- **prediction error:**

- **Notation:** define the following vectors
 - ▶ observation:

 - ▶ prediction:

 - ▶ residual:

Prediction Error: Measures

- Root-Mean-Square (RMS) prediction error:
- Relative prediction error:

Least Squares Model Fitting

- **Optimization Problem:**
- **Equivalent Least Squares Problem:** with $\hat{y}^{(i)} = f(x^{(i)})$, we write

- The j -th column of A is the j -th basis function, evaluated at each of the data points $x^{(1)}, \dots, x^{(m)}$.
- The i -th row gives the values of the p basis functions on the i -th data point $x^{(i)}$.

Least Squares Model Fitting

- Data fitting as least squares:

$$\min_{\theta} \|A\theta - \mathbf{y}\|_2^2 \quad \text{[Least Squares Problem]}$$

- **Solution:** (cf. [Mod2-L1])

- We say that the model parameter values $\hat{\theta}$ are obtained by least squares fitting on the data set.
- **Minimum square error (MSE):** $\|\mathbf{y} - A\hat{\theta}\|_2^2$
- **Minimum mean square error (MMSE):** $\frac{1}{m} \|\mathbf{y} - A\hat{\theta}\|_2^2$

Intuition Road Map

- [Univariate] Least squares fit with a constant $\rightarrow \hat{y}$ is average of data y
- [Univariate] Straight line fit (linear/affine model) $\rightarrow \hat{y}$ is linear combination of average of y and linear scaling of $(x - \text{avg}(x))$
- More complex feature maps such as polynomials, piecewise linear functions etc.
- [Multivariate] Regression

[Univariate]: Least squares fit with a constant

- **set-up:** Let $p = 1$ and $f_1(x) = 1$ for all x
- **model:** $f(x) = \theta_1$ and $A = \mathbf{1} \in \mathbb{R}^{m \times 1}$ —i.e., least squares fitting in this case is the same as choosing the best constant value θ_1 to approximate the data $y^{(1)}, \dots, y^{(m)}$
- **Solution:**
 - The RMS fit to the data (i.e., the RMS value of the optimal residual) is
 - **Interpretation:** the average value and the standard deviation of the outcomes, as the best constant fit and the associated RMS error, respectively.

[Univariate]: Straight Line Fit

- Suppose that $n = 1$ —i.e., the feature vector $x \in \mathbb{R}$ is a scalar and so is $y \in \mathbb{R}$
- **Visualization:** Data points $\{(x^{(i)}, y^{(i)})\}_{i=1}^m$ can thus be plot in the plane
- **Straight Line fit:** Let feature mappings be $f_1(x) = 1$ and $f_2(x) = x$ so that our model is

$$f(x) = \theta_1 + \theta_2 x \quad \text{[Straight-Line]}$$

- **Note:** this is often called a "linear model" however its really *affine*
- The data matrix $A \in \mathbb{R}^{m \times 2}$ is given by

[Univariate]: Straight Line Fit (P1)

- **Solution:** as before we have $\theta = (A^T A)^{-1} A^T \mathbf{y}$
- Explicit solution:

[Univariate]: Straight Line Fit (P2)

- Explicit solution:

$$\hat{\theta} = \frac{1}{m\mathbf{x}^\top\mathbf{x} - (\mathbf{1}^\top\mathbf{x})^2} \begin{bmatrix} \mathbf{x}^\top\mathbf{x} & -\mathbf{1}^\top\mathbf{x} \\ -\mathbf{1}^\top\mathbf{x} & m \end{bmatrix} \begin{bmatrix} \mathbf{1}^\top\mathbf{y} \\ \mathbf{x}^\top\mathbf{y} \end{bmatrix}$$

- Multiply by m^2/m^2 (put the m^2 in the numerator on the scalar term and the $1/m$ terms on the vector and matrix, resp.) to get

[Univariate]: Straight Line Fit

- Solution:

$$\hat{\theta} = \frac{1}{\text{rms}(\mathbf{x})^2 - \text{avg}(\mathbf{x})^2} \begin{bmatrix} \text{rms}(\mathbf{x})^2 & -\text{avg}(\mathbf{x}) \\ -\text{avg}(\mathbf{x}) & 1 \end{bmatrix} \begin{bmatrix} \text{avg}(\mathbf{y}) \\ \frac{1}{m} \mathbf{x}^\top \mathbf{y} \end{bmatrix}$$

- $\hat{\theta}_2$ is the slope and can be expressed as follows

$$\hat{\theta}_2 = \frac{m \mathbf{x}^\top \mathbf{y} - (\mathbf{1}^\top \mathbf{x})(\mathbf{1}^\top \mathbf{y})}{m \mathbf{x}^\top \mathbf{x} - (\mathbf{1}^\top \mathbf{x})^2} = \frac{(\mathbf{x} - \text{avg}(\mathbf{x})\mathbf{1})^\top (\mathbf{y} - \text{avg}(\mathbf{y})\mathbf{1})}{\|\mathbf{x} - \text{avg}(\mathbf{x})\mathbf{1}\|_2^2} = \frac{\text{std}(\mathbf{y})}{\text{std}(\mathbf{x})} \rho$$

where ρ is the correlation coefficient [VMLS, Ch. 3] defined by

$$\rho = \frac{(\mathbf{x} - \text{avg}(\mathbf{x})\mathbf{1})^\top (\mathbf{y} - \text{avg}(\mathbf{y})\mathbf{1})}{m \text{std}(\mathbf{x}) \text{std}(\mathbf{y})} \quad \text{and} \quad \text{std}(\mathbf{x}) = \frac{\|\mathbf{x} - \text{avg}(\mathbf{x})\mathbf{1}\|_2}{\sqrt{m}}$$

[Univariate]: Straight Line Fit

- $\hat{\theta}_1$ is the intercept and it also has a simple expression. \implies

$$\text{[Normal Equations]: } (A^\top A)\theta = A^\top \mathbf{y} \implies m\theta_1 + (\mathbf{1}^\top \mathbf{x})\theta_2 = \mathbf{1}^\top \mathbf{y}$$

so that

$$\hat{\theta}_1 = \text{avg}(\mathbf{y}) - \hat{\theta}_2 \text{avg}(\mathbf{x})$$

- Hence,

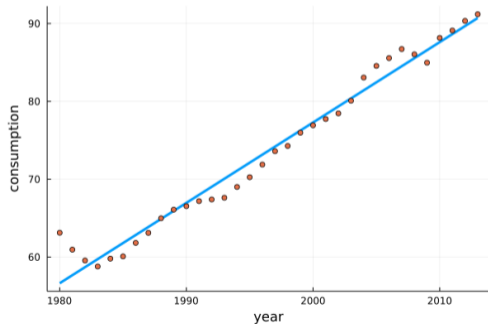
$$\hat{y} = f(x) = \text{avg}(\mathbf{y}) + \rho \frac{\text{std}(\mathbf{y})}{\text{std}(\mathbf{x})} (x - \text{avg}(\mathbf{x}))$$

- When $\text{std}(\mathbf{y}) \neq 0$ we have

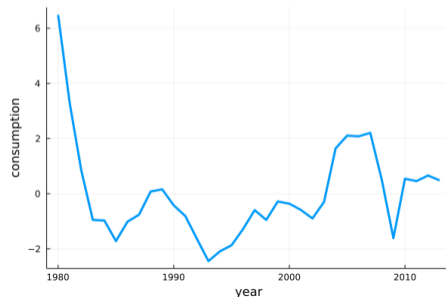
$$\frac{\hat{y} - \text{avg}(\mathbf{y})}{\text{std}(\mathbf{y})} = \rho \frac{x - \text{avg}(\mathbf{x})}{\text{std}(\mathbf{x})}$$

Example: Petroleum Consumption

- Data represents a series of samples of world petroleum consumption y at time $x^{(i)} = i$.
- straight-line fit: $\hat{y}^{(i)} = \theta_1 + \theta_2 i$, $i = 1, \dots, m$ (trend line)



- Slope: θ_2 is the *trend*
- de-trended time series: $y - \hat{y}$
- positive \rightarrow time series above straight-line fit



[Univariate]: Polynomial Fit

- A simple extension beyond the straight-line fit is a polynomial fit—i.e., where

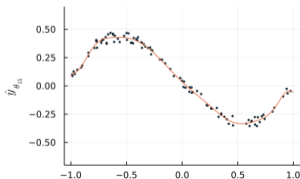
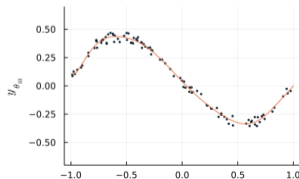
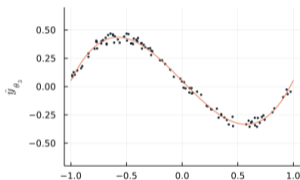
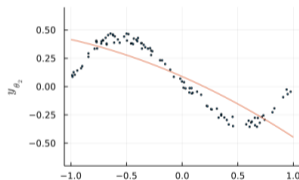
$$f_i(x) = x^{i-1}, \quad i = 1, \dots, p$$

- f is a polynomial of degree $p - 1$:

$$f(x) = \theta_1 + \theta_2 x + \dots + \theta_p x^{p-1}$$

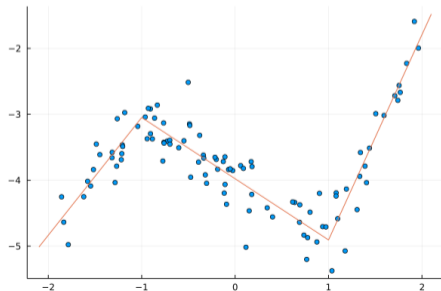
- Data matrix:

$$A = \begin{bmatrix} 1 & x^{(1)} & \dots & (x^{(1)})^{p-1} \\ 1 & x^{(2)} & \dots & (x^{(2)})^{p-1} \\ \vdots & \vdots & & \vdots \\ 1 & x^{(m)} & \dots & (x^{(m)})^{p-1} \end{bmatrix}$$



[Univariate]: Piecewise Linear Fit

- Yet another generalization is the piecewise linear fit: a piecewise linear function with knot points $c_1 < c_2 < \dots < c_k$ is a continuous function that is affine between the knot points.
- **basis functions:** $f_1(x) = 1$, $f_2(x) = x$, $f_{i+2} = (x - c_i)_+ = \max\{x - c_i, 0\}$, $i = 1, \dots, k$



[Multivariate]: Regression (Multivariate)

- Regression: $\hat{y} = x^\top \beta + v$, where β is the weight vector and v the offset [VMLS, Ch. 2]
- Let $f_1(x) = 1$ and $f_i(x) = x_{i-1}$, $i = 2, \dots, n + 1$ so that $p = n + 1$
- Transforming to data fitting model:

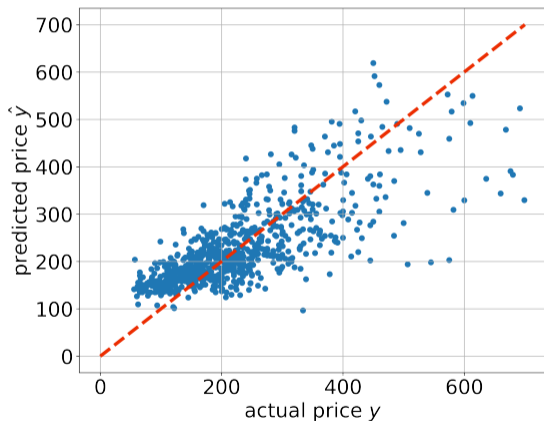
- Data matrix:

- General data fitting as regression: $\tilde{x} = (f_2(x), \dots, f_p(x))$ so that $\hat{y} = \tilde{x}^\top \beta + v$
- Comparing to linear parametric model:

$$\hat{y} = \theta_1 f_1(x) + \dots + \theta_p f_p(x) \quad \text{where } v = \theta_1, \beta = (\theta_2, \dots, \theta_p)$$

Example: House Price Prediction

- 774 house sales in Sacramento over 5 day period
- RMS fitting error is \approx \$74.8k
- Compare to standard deviation of the data which is \$112.8k
- i.e., basic regression model predicts the prices substantially better than a constant model



Validation & Generalization

[VMLS, 13.2]

Generalization

- **Goal of model fitting:** not to just achieve a good fit on the given data set, but rather to achieve a good fit on new data that we have not yet seen.
- *How well can we expect a model to predict y for future or other unknown values of x ?*
- **Generalization ability:** a model that makes reasonable predictions on new, unseen data has generalization ability, or generalizes
- **Over-fitting:** a model that makes poor predictions on new, unseen data is said to suffer from over-fit

Validation

- A simple but effective method for assessing the generalization ability of a model is called *out-of-sample validation*.
- **Method:**
 - ▶ **step 1:** divide data into training set and test (validation) set: e.g., randomly assign 80% of data pairs to training and 20% to testing
 - ▶ **step 2:** fit model only on training set
 - ▶ **step 3:** evaluate (e.g., via RMS) on test set
- **Assessment:** If the RMS prediction error on the test set is much larger than the RMS prediction error on the training set, we conclude that our model has poor generalization ability.
- **Over-fitting:** When the RMS prediction error on the training set is much smaller than the RMS prediction error on the test set, we say that the model is over-fit.

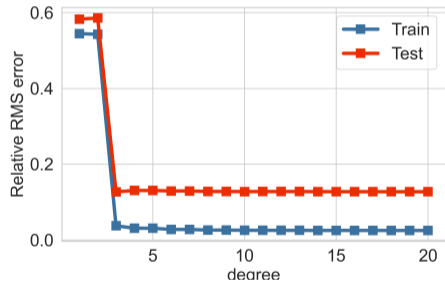
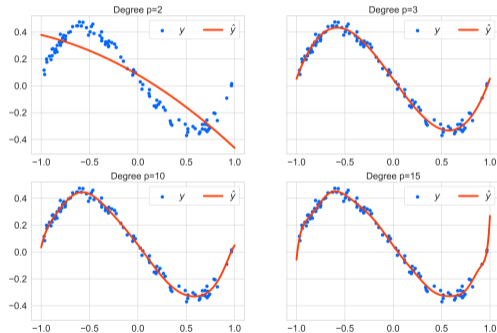
Model prediction quality & generalization ability: Not the same!

- Caution!: A model can perform poorly and yet have good generalization ability.
- e.g., let $\hat{y} = 0$ always.
- This will perform poorly on both training and test data, however the RMS will be similar meaning it has good generalization ability
- **Goal: Performance+generalization:** we seek a model that makes *good predictions on the training data set* and also makes *good predictions on the test data set*.

Generalization & Model Choice

- **Key Question:** *which is the best model?* e.g., straight-line fit, polynomial of degree 2,3,4 etc.
- **Rule of thumb 1 (least error):** we should choose a model that has test set RMS error that is near the minimum over the candidates
- **Rule of thumb 2 (simplest among them):** If multiple candidates achieve test set performance near the minimum, we should choose the 'simplest' one among these candidates.

Generalization & Model Choice: Example



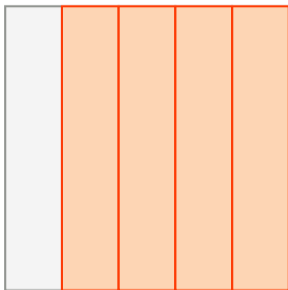
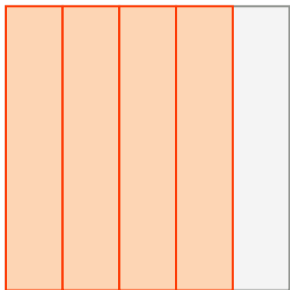
Cross-validation

- Cross-validation is an extension of out-of-sample validation that can be used to get even more confidence in the generalization ability of a model.
- **Method:** e.g., 10-fold cross-validation
 - ▶ Divide the data set into 10 sets (called folds)
 - ▶ Fit the model using folds 1–9 as training data and fold 10 as test
 - ▶ Fit the model using folds 1–8, 10 as training data and fold 9 as test
 - ▶ ⋮
 - ▶ End up with ten models and ten assessments of these models
 - ▶ RMS cross-validation error:

$$\text{rms} = \sqrt{\frac{\epsilon_1^2 + \dots + \epsilon_{10}^2}{10}}, \quad \epsilon_i := \text{rms error of model } i$$

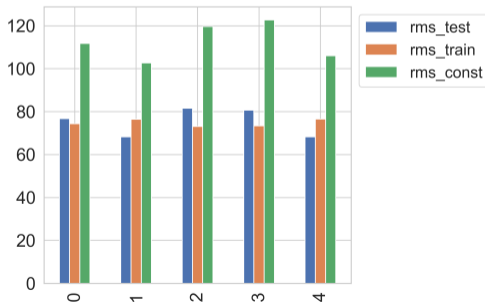
- Cross validation is used for checking of the selection of basis functions.
- Which model to use?: The models should be not too different, so the choice really should not matter much.

Cross-validation



Example: House Price Prediction

- Cross-validation to assess the generalization ability of the simple regression model of the house sales data [VMLS, §2.3]
- Current RMS error \$74.8k; cross-validation will help us answer the question of how the model might do on different, unseen houses.
- randomly partition the data set of 774 sales records into five folds, four of size 155 and one of size 154.



Limitations

- **Caution:** the basic assumption that the test data and future data are similar can (and does) fail in some applications
- e.g.: A model that predicts consumer demand, trained and validated on this year's data, can make much poorer predictions next year, simply because consumer tastes shift.
- e.g.: In finance, patterns of asset returns periodically shift, so models that predict well on test data from this year need not predict well next year.
- **Caution:** the data set is small \implies harder to interpret out-of-sample and cross-validation results.
- e.g.: out-of-sample test RMS error might be small due to good luck, or large due to bad luck, in the selection of the test set.

Feature Engineering

Feature Engineering

- (cf. [VMLS, §13.1.2]) Fitting a linear parametric model reduces to regression with new features which are the original features x mapped through the basis (or feature mapping) functions f_1, \dots, f_p .
- **Feature engineering/selection**: the process of choosing or selecting the feature maps
- To choose among candidate basis functions, we use out-of-sample validation or cross-validation.
- Heuristics:
 - ▶ Include the constant basis function $f_1(x) = 1$ (equiv. offset in basic regression model)
 - ▶ Include the original features $f_i(x) = x_{i-1}$, $i = 2, \dots, n + 1$ (equiv. starting with the basic regression model)
- **Adding more features vs. reducing them**: we can both add more features (i.e., $p > n$) to get a richer model (and select amongst them via cross-validation) or combine or reduce the number of features to reduce the dimension (i.e., $p < n$).

Transforming Features: Standardizing Features

- The first thing to consider is simple transformations of features.
- **Standardization:** scale and offset

$$f_i(x) = \frac{(x_i - b)}{a_i}, \quad i = 2, \dots, n + 1$$

- ▶ $b_i = \text{avg}(x)$
- ▶ $a_i = \text{std}(x)$
- ensures that the average of $f_i(x)$ is zero and the standard deviation is near one.
- This is called *standardizing* or *z-scoring*

Transforming Features: Winsorizing Features

- When the data include some very large values that are thought to be errors (say, in collecting the data), it is common to **clip** or **winsorize** (for Charles P. Windsor) the data
- e.g.,

$$f_i(x) = \begin{cases} x, & |x| \leq 3 \\ 3, & x > 3 \\ -3, & x < -3 \end{cases}$$

Transforming Features: Log Transform

- When feature values are positive and vary over a wide range, it is common to replace them with their logarithms.
- If the feature value also includes the value zero, then a common trick is to replace the log transform with $f_i(x) = \log(x + 1)$
- **Effect:** compresses the range of values that we encounter.

Creating New Features: Expanding Categoricals

- Some features take on only a few values, such as -1 and 1 or 0 and 1 , which might represent some value like presence or absence of some symptom (e.g., Boolean features).
 - ▶ Likert scale ([VMLS, page 71]) takes small number of values $-2, -1, 0, 1, 2$
 - ▶ Days of the week $0, 1, \dots, 6$
- **Expanding Categoricals:** replace a feature with ℓ values with a set of $\ell - 1$ new features each of which is Boolean. This simply records whether or not the original feature has the associated value.
- If all values are zero across the $\ell - 1$ new features, then the original feature had the default value
- **Example:** suppose $x_1 \in \{-1, 0, 1\}$. Using the feature value 0 as the default feature, we replace x_1 with two mapped features

$$f_1(x) = \begin{cases} 1, & x_1 = -1 \\ 0, & \text{otherwise} \end{cases} \quad f_2(x) = \begin{cases} 1, & x_1 = 1 \\ 0, & \text{otherwise} \end{cases}$$

Creating New Features: Generalized Additive model

- New features that are nonlinear functions of the original features—e.g.,

$$\min\{x_i + a, 0\} \quad \text{and} \quad \max\{x_i - b, 0\} \quad \text{where } a, b \text{ are parameters}$$

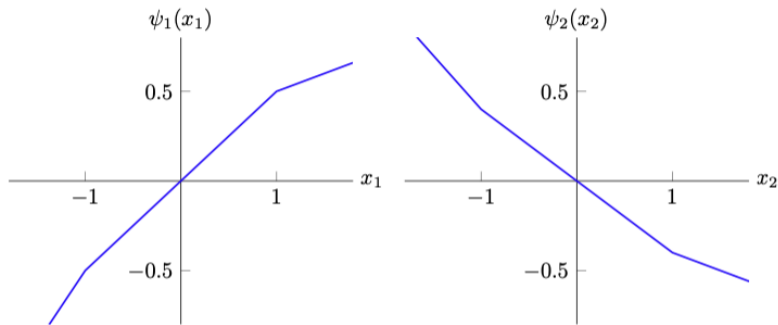
- $\min\{x_i + a, 0\}$: amount by which feature x_i is below $-a$, and
- $\max\{x_i - b, 0\}$: amount by which feature x_i is above b
- **Generalized additive model**: if the data is standardized set $a = b = 1$ so that

$$\hat{y} = \psi_1(x_1) + \cdots + \psi_n(x_n)$$

where

$$\psi_i(x_i) = \theta_{n+i} \min\{x_i + a, 0\} + \theta_i x_i + \theta_{2n+i} \max\{x_i - b, 0\}$$

Creating New Features: Generalized Additive model



- Example: $n = 2$ original features. Prediction \hat{y} is a sum of two piecewise-linear functions, each depending on one of the original features.
- $n = 2$, $a = b = 1$, and $\theta_1 = 0.5$, $\theta_2 = -0.4$, $\theta_3 = 0.3$, $\theta_4 = -0.2$, $\theta_5 = -0.3$, $\theta_6 = 0.2$.

Creating New Features: Products and Interactions

- New features can be developed from pairs of original features, for example, their product.
- e.g., $x_i x_j$ for $i, j = 1, \dots, n$ and $i \leq j$
- Product features are easily interpretable when the original features are Boolean (take the values 0 or 1)—i.e, $x_i = 1$ means that feature i is present and the new product has the value one exactly when both features i and j have occurred.

Creating New Features: Stratified Models

- **Stratified Model:** several different sub-models, and choose the one to use depending on the values of the regressors.
- e.g., instead of treating gender as a regressor in a single model of some medical outcome, we build two different sub-models, one for male patients and one for female patients.
- More generally, we can carry out clustering of the original feature vectors, and fit a separate model within each cluster.
- To evaluate \hat{y} for a new x , we first determine which cluster x is in, and then use the associated model.
- Whether or not a stratified model is a good idea is checked using out-of-sample validation.

Advanced Feature Generation

- Custom mappings: driven by application
- Predictions from other models: if there are existing models for the data we can exploit those
- Distance to cluster representatives z_1, \dots, z_k : $f_i(x) = e^{-\|x-z_i\|^2/\sigma^2}$
- Random features: random linear combination of original features
- Neural network features: computes transformed features using compositions of linear transformations interspersed with nonlinear mappings such as the absolute value.

Summary of Heuristics/Rule-of-Thumb

- **Try simple models first.** Start with a constant, then a simple regression model, and so on. You can compare more sophisticated models against these.
- **Compare competing candidate models using validation.** Adding new features will always reduce the RMS error on the training data, but the important question is whether or not it substantially reduces the RMS error on the test or validation data sets.
- **Adding new features can easily lead to over-fit.** The most straightforward way to avoid over-fit is to keep the model simple. We mention here that another approach to avoiding over-fit, called *regularization* [VMLS, ch. 15]

Example: House Price Prediction

see `Mod2-Lec2.ipynb`