

Lecture 14: Review

STATS 202: Data Mining and Analysis

Linh Tran

tranlm@stanford.edu



Department of Statistics
Stanford University

August 11, 2023



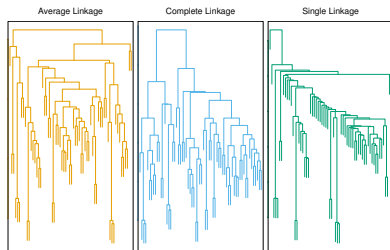
- ▶ Kaggle predictions due Sunday night.
- ▶ Final project write-up is due Wednesday.
 - ▶ **Reference your Kaggle leaderboard name on Page 1**
- ▶ Final exam is next Saturday
 - ▶ Time: Saturday August 19 @ 7:00 PM -10:00 PM
 - ▶ Location: Skilling Auditorium
 - ▶ 8 questions (lowest question dropped)
 - ▶ Practice exam to be released tonight (solutions next week)
 - ▶ Accommodation requests should already be made
- ▶ Course evaluation is up (on Canvas).



- ▶ Course review



- ▶ In unsupervised learning, all the variables are on equal standing, no such thing as an input and response.
- ▶ Clustering is typically applied
 - ▶ Hierarchical clustering (single, complete, or average linkage).
 - ▶ K -means clustering.
 - ▶ Expectation maximization (using Gaussian mixtures).



- ▶ Agglomerative algorithm produces a *dendrogram*.
- ▶ At each step we join the two clusters that are “closest”:
 - ▶ **Complete:** distance between clusters is maximal distance between any pair of points.
 - ▶ **Single:** distance between clusters is minimal distance.
 - ▶ **Average:** distance between clusters is the average distance.
- ▶ Height of a branching point = distance between clusters joined.



- ▶ The number of clusters is fixed at K .
- ▶ Goal is to minimize the average distance of a point to the average of its cluster.
- ▶ The algorithm starts from some assignment, and is guaranteed to decrease this average distance.
- ▶ This find a local minimum, not necessarily a global minimum, so we typically repeat the algorithm from many different random starting points.



We're interested in a response variable Y associated to each vector of predictors \mathbf{X} .

Regression: $f_0 = \mathbb{E}_0[Y|X_1, X_2, \dots, X_p]$

- ▶ A scalar value, i.e. $f_0 \in \mathbb{R}$
- ▶ \hat{f}_n therefore gives us estimates of y

Classification: $f_0 = \mathbb{P}_0[Y = y|X_1, X_2, \dots, X_p]$

- ▶ A vectored value, i.e.
 $f_0 = [p_1, p_2, \dots, p_K] : p_j \in [0, 1], \sum_K p_j = 1$
- ▶ n.b. In a binary setting this simplifies to a scalar, i.e.
 $f_0 = p_1 : p_1 = \mathbb{P}_0[Y = 1|X_1, X_2, \dots, X_p] \in [0, 1]$
- ▶ \hat{f}_n therefore gives us predictions of each class
- ▶ Can take the arg max, giving us Bayes Classifier



Let x_0 be a fixed point, $y_0 = f_0(x_0) + \epsilon$, and \hat{f}_n be an estimate of f_0 from $(x_i, y_i) : i = 1, 2, \dots, n$.

The MSE at x_0 can be decomposed as

$$\begin{aligned}MSE(x_0) &= \mathbb{E}_0[y_0 - \hat{f}_n(x_0)]^2 && (1) \\ &= \text{Var}(\hat{f}_n(x_0)) + \text{Bias}(\hat{f}_n(x_0))^2 + \text{Var}(\epsilon_0) && (2)\end{aligned}$$



Regression:

- ▶ MSE $((y_i - \hat{y}_i)^2)$
- ▶ AIC, BIC, R^2 , Adjusted R^2

Classification:

- ▶ Cross-entropy $((y_i \log(\hat{p}_i))$
- ▶ 0-1 loss $(\mathbb{I}(y_i \neq \hat{y}_i))$
- ▶ Confusion matrix
- ▶ Receiver operating characteristic curve (& AUC)
- ▶ Gini index $(\sum_{m=1}^{|T|} q_m \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk}))$

Misc:

- ▶ Hinge loss $(\max(0, 1 - yf))$



- ▶ Our main technique is cross-validation.
- ▶ Different approaches:
 1. **Validation set:** Split the data in two parts, train the model on one subset, and compute the test error on the other.
 2. **k -fold:** Split the data into k subsets. Average the test errors computed using each subset as a validation set.
 3. **LOOCV:** k -fold cross validation with $k = n$.
- ▶ No approach is superior to all others.
- ▶ What are the main differences? How do the bias and variance of the test error estimates compare? Which methods depend on the random seed?



- ▶ **Main idea:** If we have enough data, the empirical distribution is similar to the actual distribution of the data.
- ▶ Resampling with replacement allows us to obtain pseudo-independent datasets.
- ▶ They can be used to:
 1. Approximate the standard error of a parameter (say, β in linear regression), which is just the standard deviation of the estimate when we repeat the procedure with many independent training sets.
 2. **Bagging:** By averaging the *predictions* \hat{y} made with many independent data sets, we eliminate the variance of the predictor.

n.b. Can instead use the jackknife as a linear approximation.



- ▶ (Non-linear) feature transformations
- ▶ Standardization
- ▶ Kernels
- ▶ Neural networks
- ▶ True, empirical, estimated distributions



- ▶ Coefficients, standard errors, and hypothesis testing
- ▶ Interactions between predictors
- ▶ Non-linear relationships
- ▶ Correlation of error terms
- ▶ Non-constant variance of error (heteroskedasticity)
- ▶ Outliers
- ▶ High leverage points
- ▶ Collinearity
- ▶ Mis-specification



- ▶ Multiple linear regression
- ▶ Stepwise selection methods
- ▶ Ridge regression, Lasso, and elastic net
- ▶ Non-linear methods:
 - ▶ Polynomial regression
 - ▶ Cubic splines
 - ▶ Smoothing splines
 - ▶ Local regression
 - ▶ GAMs: Combining the above methods with multiple predictors
- ▶ Nearest neighbors regression
- ▶ Decision trees, Bagging, Random Forests, Boosting, and Neural Networks
- ▶ Neural Networks



- ▶ Nearest neighbors classification
- ▶ Naive Bayes
- ▶ Logistic regression
- ▶ LDA and QDA
- ▶ Stepwise selection methods
- ▶ Support vector classifier and support vector machines
- ▶ Decision trees, Bagging, Random Forests, Boosting
- ▶ Neural Networks



Very popular, since they give very smooth predictions over X .

- ▶ Define a set of knots $\xi_1 < \xi_2 < \dots < \xi_K$.
- ▶ We want the function $Y = f(X)$ to:
 1. Be a cubic polynomial between every pair of knots ξ_i, ξ_{i+1} .
 2. Be continuous at each knot.
 3. Have continuous first and second derivatives at each knot.

Fact: Given constraints, we need $K + 3$ basis functions:

$$f(X) = \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 + \beta_4 h(X, \xi_1) + \dots + \beta_{K+3} h(X, \xi_K) \quad (3)$$

where,

$$h(x, \xi) = \begin{cases} (x - \xi)^3 & \text{if } x > \xi \\ 0 & \text{otherwise} \end{cases}$$



Our goal is to find the function f which minimizes

$$\sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int f''(x)^2 dx$$

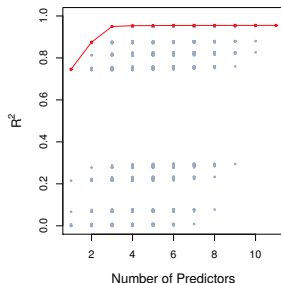
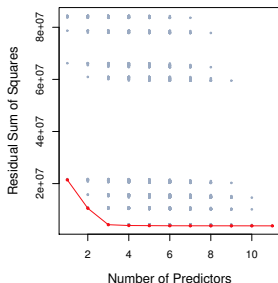
- ▶ The RSS of the model.
- ▶ A penalty for the roughness of the function.

For regularization, we have that $\lambda \in (0, \infty)$

- ▶ When $\lambda = 0$, f can be any function that interpolates the data.
- ▶ When $\lambda = \infty$, f will be the simple least squares fit



Idea: Why not just use the subset of observations *closest* to the point we're predicting at?



- ▶ Observations averaged *locally* for predictions.
- ▶ Can use different weighting kernels, e.g.

$$K_\lambda(x_0, x) = D \left(\frac{|x - x_0|}{h_\lambda(x_0)} \right) \quad (4)$$



The extension of basis functions to multiple predictors (while maintaining additivity) , e.g.

Linear model

$$\text{wage} = \beta_0 + \beta_1 \text{year} + \beta_2 \text{age} + \beta_3 \text{education} + \epsilon \quad (5)$$

Additive model

$$\text{wage} = \beta_0 + f_1(\text{year}) + f_2(\text{age}) + f_3(\text{education}) + \epsilon \quad (6)$$

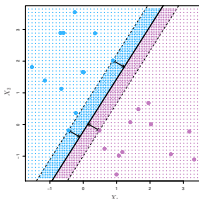
The functions f_1, \dots, f_p can be polynomials, natural splines, smoothing splines, local regressions, etc.



$$\max_{\beta_0, \beta, \epsilon} M \quad (7)$$

subject to

- ▶ $\|\beta\| = 1$
- ▶ $y_i(\beta_0 + x_i^T \beta) \geq M(1 - \epsilon_i) \quad \forall i = 1, \dots, n$
- ▶ $\epsilon_i \geq 0 \quad \forall i = 1, \dots, n$ and $\sum_{i=1}^n \epsilon_i \leq C$



n.b. Can use kernels to capture non-linearities.



Using a *greedy* approach:

- ▶ Start with a single region R_1 , and iterate:
 - ▶ Select a region R_k , a predictor X_j , and a splitting point s , such that splitting R_k with the criterion $X_j < s$ produces the largest decrease in RSS:

$$\sum_{m=1}^{|T|} \sum_{x_i \in R_m} (y_i - \bar{y}_{R_m})^2$$

- ▶ Redefine the regions with this additional split.
- ▶ Terminate when there are 5 observations or fewer in each region.
- ▶ This grows the tree from the root towards the leaves.



Algorithm 15.1 *Random Forest for Regression or Classification.*

1. For $b = 1$ to B :
 - (a) Draw a bootstrap sample \mathbf{Z}^* of size N from the training data.
 - (b) Grow a random-forest tree T_b to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached.
 - i. Select m variables at random from the p variables.
 - ii. Pick the best variable/split-point among the m .
 - iii. Split the node into two daughter nodes.
2. Output the ensemble of trees $\{T_b\}_1^B$.

To make a prediction at a new point x :

Regression: $\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$.

Classification: Let $\hat{C}_b(x)$ be the class prediction of the b th random-forest tree. Then $\hat{C}_{\text{rf}}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$.



Boosting uses a set of weak learners (e.g. decision trees) to create a strong one.

The general algorithm is:

1. Fit an initial \hat{f}_n^0 to the data and compute residuals r_i .



Boosting uses a set of weak learners (e.g. decision trees) to create a strong one.

The general algorithm is:

1. Fit an initial \hat{f}_n^0 to the data and compute residuals r_i .
2. For $b = 1, \dots, B$:
 - ▶ Fit a weak learner \hat{f}_n^b on the residuals.



Boosting uses a set of weak learners (e.g. decision trees) to create a strong one.

The general algorithm is:

1. Fit an initial \hat{f}_n^0 to the data and compute residuals r_i .
2. For $b = 1, \dots, B$:
 - ▶ Fit a weak learner \hat{f}_n^b on the residuals.
 - ▶ With learning rate λ_b , update prediction to:

$$\hat{f}_n \leftarrow \hat{f}_n + \lambda_b \hat{f}_n^b. \quad (8)$$



Boosting uses a set of weak learners (e.g. decision trees) to create a strong one.

The general algorithm is:

1. Fit an initial \hat{f}_n^0 to the data and compute residuals r_i .
2. For $b = 1, \dots, B$:

- ▶ Fit a weak learner \hat{f}_n^b on the residuals.
- ▶ With learning rate λ_b , update prediction to:

$$\hat{f}_n \leftarrow \hat{f}_n + \lambda_b \hat{f}_n^b. \quad (8)$$

- ▶ Update the residuals

$$r_i \leftarrow r_i - \lambda_b \hat{f}_n^b(x_i). \quad (9)$$



Boosting uses a set of weak learners (e.g. decision trees) to create a strong one.

The general algorithm is:

1. Fit an initial \hat{f}_n^0 to the data and compute residuals r_i .
2. For $b = 1, \dots, B$:
 - ▶ Fit a weak learner \hat{f}_n^b on the residuals.

- ▶ With learning rate λ_b , update prediction to:

$$\hat{f}_n \leftarrow \hat{f}_n + \lambda_b \hat{f}_n^b. \quad (8)$$

- ▶ Update the residuals

3. Output prediction, e.g. $r_i \leftarrow r_i - \lambda_b \hat{f}_n^b(x_i).$ (9)

$$\hat{f}_n(x) = \hat{f}_n^0 + \sum_{b=1}^B \lambda_b \hat{f}_n^b(x). \quad (10)$$



Neural networks are simply a generalization of the logistic regression case, e.g. for

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2) \quad (11)$$



Neural networks are simply a generalization of the logistic regression case, e.g. for

$$\mathbb{P}(Y = 1|\mathbf{X}) = \sigma(\sigma(\mathbf{X}\mathbf{W}_1)\mathbf{W}_2) \quad (11)$$

Our loss is

$$L(y_i, f(\mathbf{X}_i)) = -y_i \log(p_i) - (1 - y_i) \log(1 - p_i), \text{ where} \quad (12)$$

$$p_i = \frac{1}{1 + \exp(-Z_{2,i})} \quad (13)$$

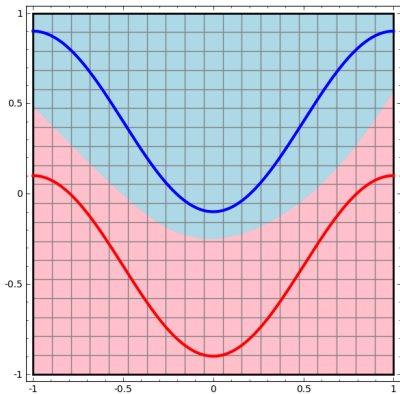
$$Z_{2,i} = h_i \mathbf{W}_2 \quad (14)$$

$$h_i = \frac{1}{1 + \exp(-Z_{1,i})} \quad (15)$$

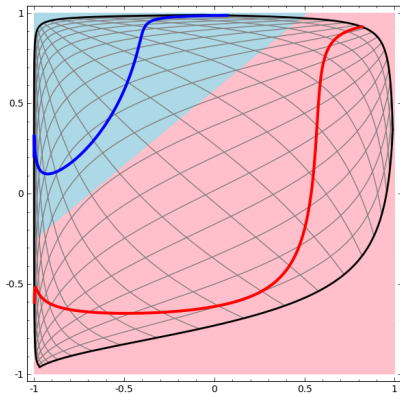
$$Z_{1,i} = \mathbf{X}\mathbf{W}_1 \quad (16)$$



How do the feature transformations get learned?



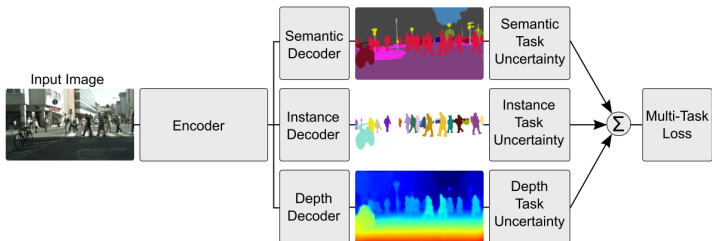
Original representation of curves



Hidden layer representation of curves

Well demonstrated by *Chris Olah's blog*.

Neural networks can be applied over multiple tasks (i.e. multi-task learning), e.g.



Kendall et al. 2017's multi-task model



Analyzing right censored survival time

- ▶ Our observed time is $Y = \min(T, C)$
- ▶ We have an associated indicator $\delta = \mathbb{I}(T \leq C)$

Two commonly used estimators

- ▶ Kaplan Meier Estimator: estimates the survival function for a small number of groups
 - ▶ Can use log-rank test to confirm statistical significance.
- ▶ Cox-proportional hazards: assumes proportionality in the hazard functions.
 - ▶ Similar to (pooled) logistic regression (breaking follow-up time into individual time ranges)



For each of the regression and classification methods:

1. What are we trying to optimize?
2. What does the fitting algorithm consist of, roughly?
3. What are the tuning parameters, if any?
4. How is the method related to other methods, mathematically and in terms of bias, variance?
5. How does rescaling or transforming the variables affect the method?
6. In what situations does this method work well? What are its limitations?