

Machine Learning I

Lecture 14: Dimensional Reduction

Nathaniel Bade

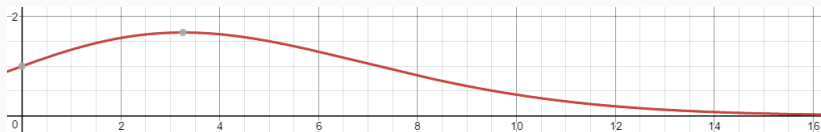
Northeastern University Department of Mathematics

Table of contents

1. Dimensional Reduction
2. Factor Analysis
3. Principle Component Analysis
4. Nonlinear PCA

Dimensional Reduction

Computational Concerns



Volume of Unit Ball

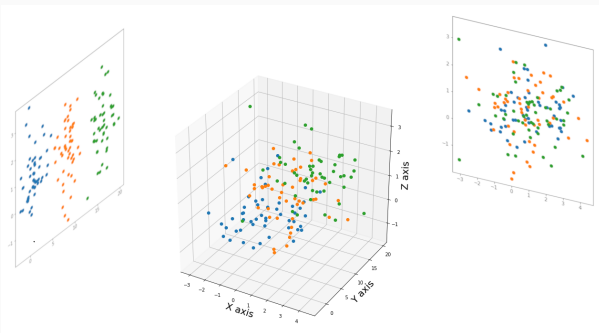
Three facts to keep in mind about the curse of dimensionality:

Volume of the unit sphere in even dimensions is $V_{2k} = \frac{\pi^k}{k!}$.

In a 1×1 square, a point has a 0.4% chance of being within .001 of the boundary. In the 10,000 dimensional unit cube the probability is greater than 99.999999%.

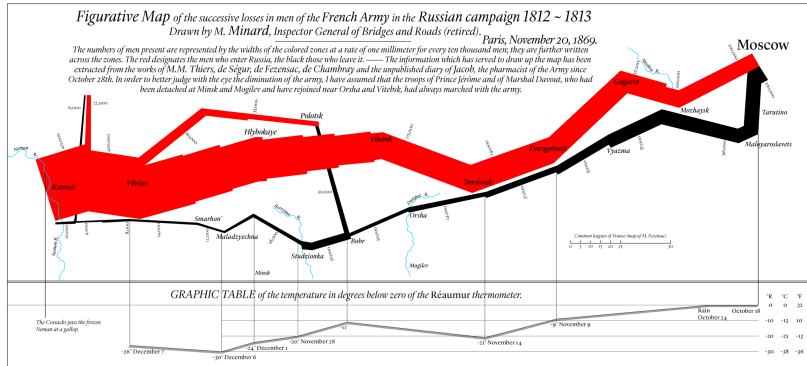
A grid search with a division at $1/2$ in each dimension contains $2^{10,000}$ regions.

Low Dimensional Projection and Meaning



In addition, low dimensional projections are *meaningful*. In fact, all of data visualization can be said to be the result careful low dimensional projections. One goal is to discover the small set of variables that controls the larger phenomena.

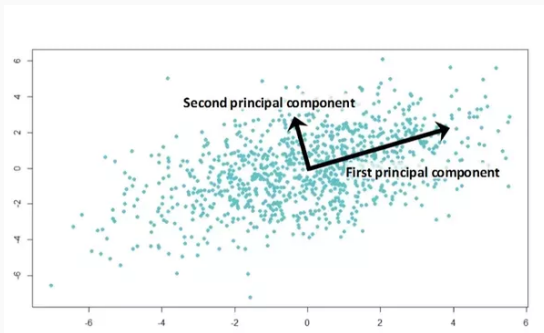
Low Dimensional Projection and Meaning



For example, Charles Minard's map of Napoleon's march into Russia is said to be one of the best examples of dimensional reduction, with 8 variables being represented using 2 dimensions, and effectually along 1.

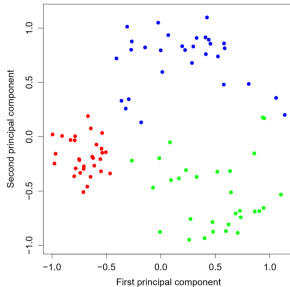
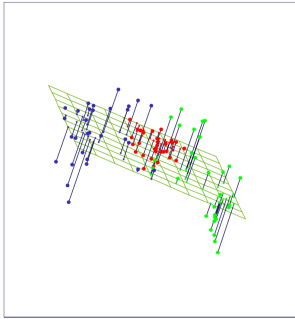
https://en.wikipedia.org/wiki/Charles_Joseph_Minard

Projections



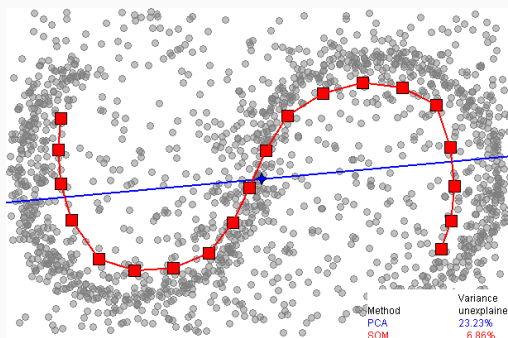
There are two main techniques in dimensional reduction: **projection onto linear subspaces** and **manifold learning**. In projection onto linear subspaces, we try to discover the linear combination of features that provide the clearest coordinate system for the dataset, i.e. the **component factors**.

Projections



In an alternative view, factor analysis tries to fit a linear subspace to a data set in such a way as to maximize the **projected variance**, that is maximize the amount of data variance captured in orthogonal project to the new space.

Projections



In **manifold learning**, we try to fit a more complicated manifold to the underlying data. Manifold learning is of course much more complicated than simple linear dimensional reduction, but here do exist good algorithms which we will talk more about in a week or two.

Factor Analysis

The Formal Problem

Given N data points $x_i \in \mathbb{R}^p$, we want to find a linear transform U that reduces x_i from a p dimensional feature space to a $k < p$ dimensional feature space without distorting the relationship between points.

Practically, this means we want to preserve the inner product:

$$\langle Ux_i, Ux_j \rangle = \langle x_i, x_j \rangle .$$

Therefore U is a $k \times p$ isometry. For the Cartesian inner product the isometries are the orthogonal matrices $U^T U = I$. Note that $UU^T \neq I$, since this would imply that lowering dimension lost no information.

The Formal Problem

Given N data points $x_i \in \mathbb{R}^p$, we want to find a linear transform U that reduces x_i from a p dimensional feature space to a $k < p$ dimensional feature space without distorting the relationship between points.

Practically, this means we want to preserve the inner product:

$$\langle Ux_i, Ux_j \rangle = \langle x_i, x_j \rangle.$$

Therefore U is a $k \times p$ isometry. For the Cartesian inner product the isometries are the orthogonal matrices $UU^T = I_k$. Note that $U^T U \neq I_p$, since we are strictly losing information when we project to a lower dimensional subspace.

The Formal Problem

The goal then is to find a new set of vectors $z_i = Ux_i$ that best represent the data, that is find the k directions in \mathbb{R}^p the best represent x_i .

We will proceed by optimizing two objectives:

Minimize the **Reconstruction error**.

Maximize the **projected variance**.

We will see that these are in fact equivalent conditions.

The Reconstruction Error

We can think about $z_i = Ux_i$ as an **encoding** or **compression** of x_i . Undoing the linear transformation then amounts to a decoding

$$\tilde{x}_i = U^T z_i = U^T U x_i .$$

The **Reconstruction error** is the difference between the original data and the decoding:

$$\sum_{i=1}^N \|x_i - \tilde{x}_i\| = \sum_{i=1}^N \|x_i - U^T U x_i\| ,$$

and the problem can be stated as finding

$$\min_{U \in \mathbb{R}^{k \times p}} \sum_{i=1}^N \|x_i - U^T U x_i\|^2 , \quad UU^T = I_k .$$

The Projected Variance

Similarly, we may try to minimize the **projected sample variance** $\widehat{\text{var}}[Ux]$. The sample variance is given as usual by

$$\widehat{\text{var}}[Ux] = \widehat{E}[(Ux)^2] - \widehat{E}[Ux]^2.$$

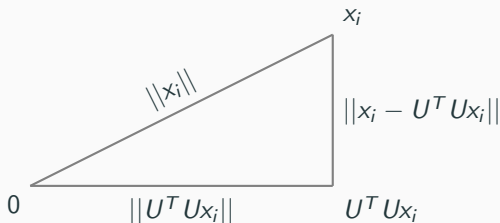
Assume our training data has been centered, that is $\widehat{E}[(x_i)_j] = 0$ for each feature $j = 1, \dots, p$. Then the sample variance is

$$\widehat{\text{var}}[Ux] = \widehat{E}[(Ux)^2].$$

Our goal then is to find

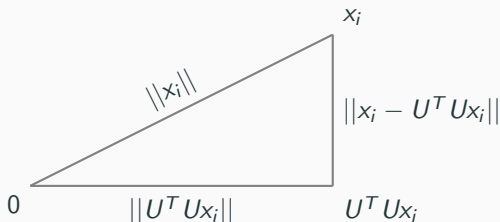
$$\max_{U \in \mathbb{R}^{k \times p}} \widehat{E}[(Ux)^2], \quad UU^T = I_k.$$

Equivalence of Conditions



The diagram above shows that to see that the two conditions are equivalent, we will show that the variance in the data (which is fixed) is the sum of the **projected sample variance** and the **reconstruction error**.

Equivalence of Conditions



Taking the expectation of the Pythagorean decomposition

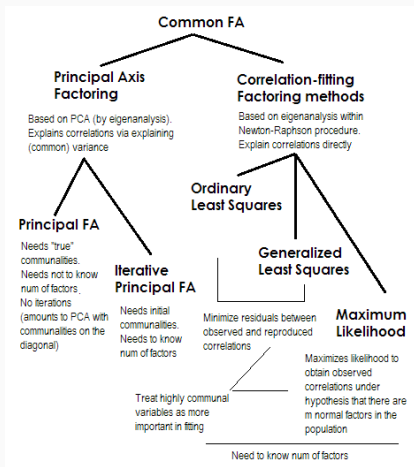
$$x_i = U^T U x_i + x_i - U^T U x_i$$

yields

$$\underbrace{\hat{E}[\|x_i\|^2]}_{\text{Sample Variance}} = \underbrace{\hat{E}[\|U^T U x_i\|^2]}_{\text{Projected Variance}} + \underbrace{\hat{E}[\|x_i - U^T U x_i\|^2]}_{\text{Reconstruction Error}}$$

We see the minimizing reconstruction error is equivalent to maximizing projected variance.

Techniques in Factor Analysis



There is no consensus on how to perform factor analysis. The problem mainly comes down to choosing the how to determine the number of factors to minimize on.

<https://stats.stackexchange.com/questions/50745/>

Principle Component Analysis

Principle Component Analysis

The most common technique for factor analysis is **Principle Component Analysis** or **PCA**. It is in fact so common that "factor analysis" usually means "factor analysis that isn't PCA."

PCA is a greedy algorithm with a beautiful mathematical interpretation. The idea is to proceed iteratively:

- Find principle component (direction) that accounts for the largest possible variance.

- Project onto the subspace orthogonal to that vector.

- Repeat, collecting the orthogonal vectors into the rows of U until U is a $p \times p$ orthogonal matrix.

Principle Component Analysis

Let x_i be demeaned data points and let \mathbf{X} be the matrix whose rows are x_i^T . We want to find a vector $w_{(1)} = (w_1, \dots, w_p)_{(1)}$ with $\|w_{(1)}\| = 1$ that maximizes the projected variance. Projection of x_i onto $w_{(1)}$ is given by

$$\text{proj}_{w_{(1)}}(x_i) = \frac{\langle x_i, w_{(1)} \rangle}{\|w_{(1)}\|^2} w_{(1)} = \langle x_i, w_{(1)} \rangle w_{(1)}.$$

Since x_i have mean 0, the variance along $w_{(1)}$ is just

$$\hat{E}[\|\langle x_i, w_{(1)} \rangle w_{(1)}\|^2] = \hat{E}[\langle x_i, w_{(1)} \rangle^2].$$

Writing $\langle x_i, w_{(1)} \rangle = x_i^T w_{(1)}$ for the Cartesian inner product, our goal is to find

$$w_{(1)} = \operatorname{argmax}_{\|w_{(1)}\|=1} \left(\sum_{i=1}^N (x_i \cdot w_{(1)})^2 \right) = \operatorname{argmax}_{\|w_{(1)}\|=1} (\|\mathbf{X} w_{(1)}\|^2).$$

Principle Component Analysis

Writing the Euclidean norm as matrix multiplication, we see that it is equivalent to find

$$w_{(1)} = \operatorname{argmax}_{\|w_{(1)}\|=1} (\| \mathbf{X} w_{(1)} \|^2) = \operatorname{argmax}_{\|w_{(1)}\|=1} (w_{(1)}^T \mathbf{X}^T \mathbf{X} w_{(1)}) .$$

The maximum is not so clear using our usual matrix differentiation. Indeed, the condition that $\|w_{(1)}\| = 1$ plays an important role in the expression even having a maximum. We can use a Lagrangian procedure to rewrite the above as a single maximization problem:

$$\mathcal{L} = w_{(1)}^T \mathbf{X}^T \mathbf{X} w_{(1)} + \lambda(1 - w_{(1)}^T w_{(1)}) .$$

(Homework) Show that the stationary points of \mathcal{L} are the solutions to the eigenvalue equation $\mathbf{X}^T \mathbf{X} w_{(1)} = \lambda w_{(1)}$.

Principle Component Analysis

Since the unit eigenvectors of $\mathbf{X}^T \mathbf{X}$ extremize $w_{(1)}^T \mathbf{X}^T \mathbf{X} w_{(1)}$,

$$\operatorname{argmax}_{\|w_{(1)}\|=1} (w_{(1)}^T \mathbf{X}^T \mathbf{X} w_{(1)})$$

is the eigenvector of $\mathbf{X}^T \mathbf{X}$ with the largest eigenvalue λ . Furthermore,

$$\widehat{\operatorname{var}}(\mathbf{X} x_{(1)}) = \frac{1}{N} w_{(1)}^T \mathbf{X}^T \mathbf{X} w_{(1)} = \frac{1}{N} w_{(1)}^T (\lambda w_{(1)}) = \frac{\lambda}{N}.$$

But recall that the covariance matrix of the data

$$\widehat{\operatorname{Cov}}(\mathbf{X}) = \frac{1}{N} \mathbf{X}^T \mathbf{X}.$$

We have shown that the first principle component is the eigenvector of the covariance matrix with the highest eigenvalue.

Principle Component Analysis

The result extends to the subsequent principle components. The proof follow exactly ass for the first principle component and in fact constructs the orthogonal decomposition

$$\widehat{\text{Cov}}(\mathbf{X}) = W\Lambda^{(p)}W^T,$$

where $\lambda^{(p)} = \text{diag}(\lambda_1, \dots, \lambda_p)$, with $\lambda_i > \lambda_{i+1}$.

Computationally, the complexity is $O(Np^2 + p^3)$, the sum computing the covariance and eigenbasis respectively.

Truncated Principle Component Analysis

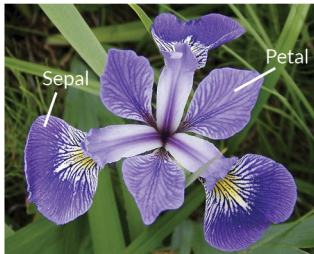
For dimensional reduction, we may stop this process at any step, revealing a **truncated** PCA matrix. Keeping only the first k principle components yields a orthogonal matrix W_k which projects x_i to the space spanned by the first k principle vectors.

Mathematically, this performs decomposition of the covariance matrix into

$$\widehat{\text{Cov}}(\mathbf{X}) = W\Lambda^{(k)}W^T,$$

where $\Lambda^{(k)} = \text{diag}(\lambda_1, \dots, \lambda_k)$, with $\lambda_i > \lambda_{i+1}$.

Machine Learning Examples: Iris Classification



Iris Versicolor



Iris Setosa



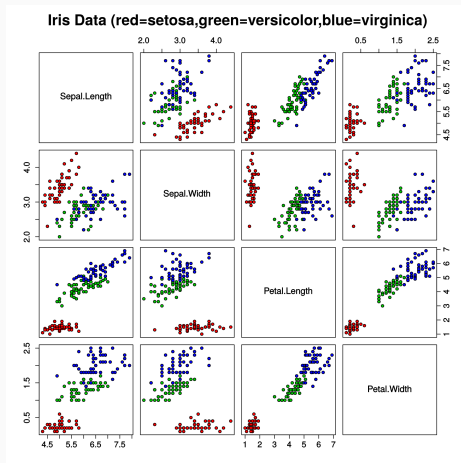
Iris Virginica

Sepal Length	Sepal Width	Petal Length	Petal Width	Species
5.1	3.5	1.4	0.2	I. Setosa
7.0	3.2	4.7	1.4	I. Versicolor

⋮

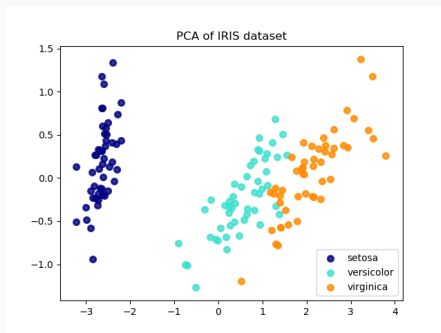
Take as an example projection onto the first two principle component in the Iris data set from Lecture 1.

Example: PCA for Iris



There are four variables, with correlations given above.

Example: PCA for Iris

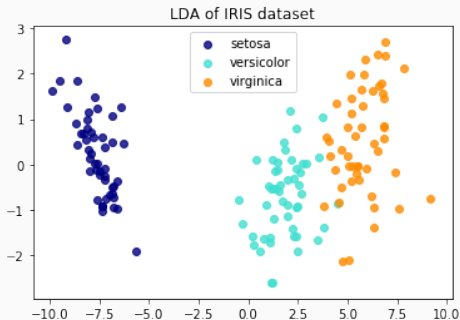


Using PCA to project onto the first two principle components yields a projection onto

$$w_{(1)} = (0.36, -0.08, 0.86, 0.36), \quad w_{(2)} = (0.66, 0.73, -0.18, -0.07)$$

We see the species are separable knowing only their features, and furthermore we have concrete measurement ratios to determine species.

Example: PCA for Iris



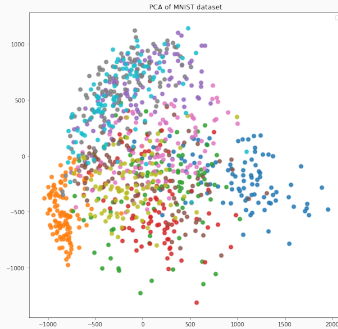
Compare for a moment with the PCA projection. LDA requires knowing the labels and requires significantly more computational time for a similar fit.

Example: PCA MNIST



The MNIST data set may provide one of the most startling examples. A simple PCA can be computed quite quickly and yields a wealth of information about the structure of the data set.

Example: PCA MNIST



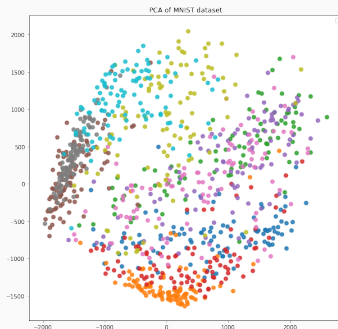
The MNIST data set may provide one of the most startling examples. A simple PCA can be computed quite quickly and yields a wealth of information about the structure of the data set. Projecting onto the first two components shows a lot of structure on the MNIST dataset.

Example: PCA MNIST



This is not always as useful. Performing PCA on the fashion MNIST data set doesn't yield as strong of forms as on the digits MNIST.

Example: PCA MNIST

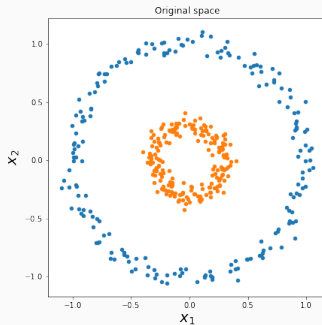


This is not always as useful. Performing PCA on the fashion MNIST data set doesn't yield as strong of forms as on the digits MNIST. On the other hand, we see much more structure and organization in the projection onto eigenvectors.

For an excellent article on visualization of MNIST, see <http://colah.github.io/posts/2014-10-Visualizing-MNIST/>

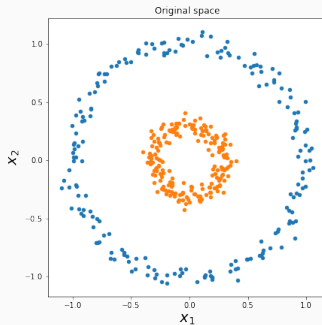
Nonlinear PCA

Nonlinear PCA



Of course, linear methods may be almost useless if the underlying structure isn't linear. For truly unknown structure we must construct smoothing maps using splines, random graphs, or other high complexity techniques. For example, in the distribution above linear PCA will do nothing to reduce the complexity of the data set.

Nonlinear PCA



However, whenever we have linear methods there is the hope that we can extend them in a computationally efficient manner using the kernel trick.

The kernel trick was to notice that if our optimization only refers to the higher dimensional feature space by way of the inner product, then we can often replace it with a much more computationally efficient kernel $k(x_i, x_j)$.

In keeping with the literature, let $\phi : \mathcal{X} \rightarrow \tilde{\mathcal{X}}$ be an embedding of the feature space into a higher dimensional space (say, the space of degree d polynomials in the features). Assume that projection has 0 mean.

The covariance matrix in the embedded space can be written

$$\tilde{C} = \widetilde{\text{Cov}} = \frac{1}{N} \sum_{i=1}^N \phi(x_i) \phi(x_i)^T,$$

with eigenvectors w_k and eigenvalues $\tilde{C} w_k = \lambda_k w_k$.

Since

$$\tilde{C} w_k = \frac{1}{N} \sum_{i=1}^N \phi(x_i) (\phi(x_i)^T w_k) = \lambda_k w_k,$$

we can write

$$w_k = \sum_{i=1}^N a_{ki} \phi(x_i).$$

If there exists $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$, then for any $\phi(x_\ell)^T$ the first equation can be written

$$\phi(x_\ell)^T \frac{1}{N} \sum_{i=1}^N \phi(x_i) (\phi(x_i)^T \sum_{j=1}^N a_{kj} \phi(x_j)) = \lambda_k \phi(x_\ell)^T \sum_{j=1}^N a_{kj} \phi(x_j),$$

or

$$\frac{1}{N} \sum_{i=1}^N K(x_\ell, x_i) \sum_{j=1}^N a_{kj} K(x_i, x_j) = \lambda_k \sum_{j=1}^N a_{kj} K(x_\ell, x_j).$$

Writing $\mathbf{K}_{ij} = K(x_i, x_j)$, and $a_k = [a_{k1}, \dots, a_{kN}]^T$

$$\frac{1}{N} \sum_{i=1}^N K(x_\ell, x_i) \sum_{i=1}^N K(x_i, x_j) = \lambda_k \sum_{j=1}^N K(x_\ell, x_j),$$

can be written

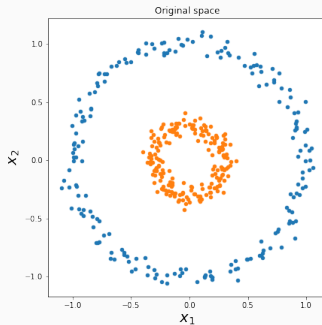
$$\mathbf{K}^2 a_k = \lambda_k N \mathbf{K} a_k.$$

So the higher dimensional fitting is solved by the kernel eigenvalue problem

$$\mathbf{K} a_k = \lambda_k N a_k.$$

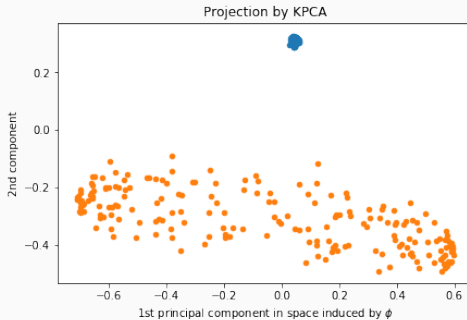
Again, this trick is restricted to transformations for which there exists a computable kernel, but if one exists this is a remarkable reduction in complexity.

Nonlinear PCA



Applying the radial basis function kernel PCA to the concentric circle data set results in a dramatic clustering.

Nonlinear PCA



Applying the radial basis function kernel PCA to the concentric circle data set results in a dramatic clustering.

References

References for this lecture: BS, Chapters 15, 16, 18 and ESII, Chapters 9, 12.

For more about SVM uses and generalization, see ESII Chapter 12. For more about SVM proofs and complexity computations, see SB Chapter 26.

More hands on with random forests in R: <https://cran.r-project.org/web/packages/rpart/vignettes/longintro.pdf>

Image reference: Comparison of Tree Loss functions:
<https://github.com/rasbt/python-machine-learning-book>