

Document Template

Template for the kaobook Class

Johnny B. Goode

October 4, 2023

An Awesome Publisher

Abstract

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

Contents

Contents	iii
1. Ensemble Learning	1
2. The Ensemble Generalisation Error in terms of Diversity	2
2.1. Bias, Variance and their Effects	2
2.1.1. Deriving Bias- and Variance-Effect	2
2.2. Diversity and Diversity-Effect	4
2.2.1. Measures of ensemble diversity	4
2.2.2. Ambiguity	4
2.2.3. A Bias-Variance-Decomposition for Ensembles	5
2.2.4. Bregman Divergences	5
3. Random Forests	7
A. Appendix	8
APPENDIX	9
B. Some more blindtext	10

1. Introduction

2. Ensemble Learning

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

2.1. Ensemble Learning Methods

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

2.2. Applications

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

2.3. Advantages of Ensemble Methods

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

3. Random Forests

This is some chapter overview. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

3.1. Decision Trees

3.1.1. Decision Trees greedily minimise loss functions

3.2. The Random Forest scheme

3.2.1. Bootstrapping & Bagging

3.2.2. Feature & Split selection

3.2.3. Number of trees

3.2.4. Depth of trees

3.2.5. Random Forests converge

3.2.6. Random Forests do not overfit

3.2.7. Random Forests are consistent

4. The Ensemble Generalisation Error in terms of Diversity

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

We begin this chapter by deriving the classical Bias-Variance-decomposition from first principles. ... Understanding the nature and the generality of bias- and variance-*effect* will allow us to very naturally and simply develop a notion of ensemble diversity as a measure of model fit, just like bias and variance.

4.1. Preliminaries

4.1.1. Other Preliminaries

Random variables, expectations, expectations are approximated by means etc.

4.1.2. Supervised Learning

Our goal is to find an algorithm that is able to map objects from \mathcal{X} to outcomes in \mathcal{Y} . Objects are described by their *features*. These are commonly numerical and we have several of these, so \mathcal{X} can be thought of as \mathbb{R}^d where d is the number of features. We will call such a representation of an object an *example*. In *classification* problems, the outcomes are discrete among k possible outcomes and we refer to them as *labels* or *classes*. For sake of simplicity, we identify these with integers, i.e. \mathcal{Y} can be thought of as the set $\{1, \dots, k\}$. In *regression* problems, the outcomes are continuous and we refer to them as *estimates*. We can think of \mathcal{Y} as \mathbb{R} . The desired mapping $q : \mathcal{X} \rightarrow \mathcal{Y}$ may be nontrivial such that it is not feasible to come up with explicit rules of how to map examples to outcomes. However, we may try to infer such rules from a given set of examples and their known outcomes. More specifically, we want to find a *learning algorithm* $\mathcal{A}(D)$ that is able to produce a set of such rules – referred to as a *model* q – given an arbitrary *training dataset* D of examples and outcomes. This is known as *supervised learning*, which we will be considering herein. To be useful, the produced model should not only be able to accurately estimate the outcomes for the given training examples, but also provide reasonable predictions for new, previously unseen examples.

To analyse the problem, we will presuppose a probability distribution $P(\mathcal{X}, \mathcal{Y})$ from which realisations of example-training pairs are drawn. We write this as $(X, Y) \sim P(\mathcal{X}, \mathcal{Y})$. This distribution is unknown in practise – else the problem would be trivially solved already. In order for our solution to be widely applicable, we will strive to make as few assumptions about P as possible. Our given dataset $\{(x_i, y_i)\}_{i=1}^n$ can be

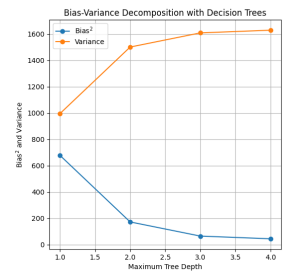


Figure 4.1.: foo!

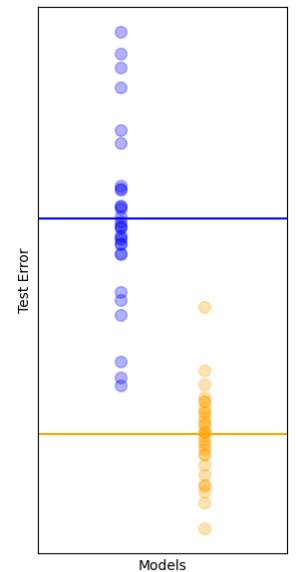


Figure 4.2.: Visualising the variance of **Decision Tree** and **Random Forest** models. Each glyph corresponds to the test error of one model trained on a random subset of the full available data. The variation of the test error around the mean test error across many dataset samples is exactly the variance. Not only do Random Forests show lower test errors on average, they seem to also have lower variance. We will explain this observation in ??

considered a random variable D drawn from $P(\mathcal{X}, \mathcal{Y})^n$ where n is the number of data points. The model q produced by our learning algorithm \mathcal{A} is a function $\mathcal{X} \rightarrow \mathcal{Y}$ that depends on the given dataset D . The prediction $q_D(X)$ is a random variable that depends on the random variables D and X .¹

1: (todo) also other sources of randomness such as random initialisation etc

As noted previously, we want to consider how well the expected classifier generalises to unseen examples, i.e. how well it performs on newly drawn example-label pairs $(X, Y) \sim P$. Such pairs can be thought of as "testing" pairs. We assess the quality of a single prediction with a *loss function* $L : \mathcal{Y} \rightarrow \mathcal{Y}$ whose value should be low if the predicted outcome is close to the true outcome. The *risk* of model q_D is the expected loss of the model over all testing points:

$$\text{Risk}(q_D) =_{\text{def}} \mathbb{E}_{(X,Y) \sim P} [L(Y, q_D(X))]$$

Since D is just a sample drawn from P , we are ultimately interested in the *expected risk*, also known as *generalisation error* or *test error*:

$$\text{GE}(\mathcal{A}) =_{\text{def}} \mathbb{E}_D [\text{Risk}(q_D)] = \mathbb{E}_{(X,Y), D} [L(Y, q_D(X))]$$

In the following, we will analyse this quantity and how and why some specific learning algorithms minimise it.

4.2. Bias, Variance and their Effects

4.2.1. Deriving Bias- and Variance-Effect

We will begin by considering the widely known bias-variance-decomposition for regression using the squared-error loss $L(y, y') = (y - y')^2$. For sake of clarity, we will from now on take the dependence of $q_D(X)$ on X as understood and write only q_D .

The variance of q_D with respect to the squared-error loss is commonly defined as the statistical variance of the regressor around its expected value. This is the expected squared error between a random value (q_D) and the closest non-random value ($\mathbb{E}_D [q_D]$).

$$\text{Var}(q_D) := \mathbb{E}_D [(q_D - \mathbb{E}_D [q_D])^2] = \min_z \mathbb{E}_D [(q_D - z)^2]$$

As such, $\mathbb{E}_D [q_D]$ is a centroid to the different realisations of q_D with respect to the loss function $L(y, x) = (y - x)^2$. This non-random centroid will turn out to be particularly interesting. For a random variable Z , we define

Definition 4.2.1 The *systematic part* of a random variable Z with respect to a loss function L is defined as

$$Z^* := \arg \min_z \mathbb{E} [L(Y, z)]$$

where the expectation is over the distribution of Z .

Note that

- ▶ $y^*(X) = \arg \min_z \mathbb{E}_Y [L(Y, z_Y(X))] = \mathbb{E}_{Y|X} [Y]$ is the *expected label*
- ▶ $q^*(X) = \arg \min_z \mathbb{E}_D [L(q_D(X), z_D(X))] is the *expected model*.$

As such, for the squared-error loss, we can write

$$\text{Var}(q_D) = \mathbb{E}_D [L(q_D, q^*)]$$

The well-known bias-variance decomposition for the squared-error loss is given as follows.

$$\begin{aligned}\mathbb{E}_{(X,Y),D} [(Y - q_D(X))^2] &= \mathbb{E}_{(X,Y)} [(Y - y^*(X))^2] + \mathbb{E}_{X,D} [(y^*(X) - q_D(X))^2] \\ &= \underbrace{\mathbb{E}_{(X,Y)} [(Y - y^*(X))^2]}_{\text{Var}(Y)} + \underbrace{\mathbb{E}_{X,D} [(q^*(X) - q_D(X))^2]}_{\text{Var}(q)} + \underbrace{\mathbb{E}_X [(y^*(X) - q^*(X))^2]}_{\text{Bias}^2(Y,q)}\end{aligned}$$

This decomposition is usually derived by expanding the square **[todo]**. The cross-terms then vanish due to that $\mathbb{E}_D [q_D] - q^* = 0$ and $y^* - \mathbb{E}_{Y|X} [Y] = 0$. This is but a special case of a more general structure applying to a certain class of losses. We will provide a more general proof later.

The first term, $\text{Var}(Y)$ is independent of D and q_D . This means we have no means of influencing it with our choice of q_D . $\text{Var}(Y)$ is also referred to as *noise*, *bayes error* or *irreducible error*. The second term, $\text{Var}(q_D)$ measures the variance of our model around its non-random centroid with respect to different realisations of the random training dataset D . This can be understood as a measure of spread of the learning algorithm with respect to different realisations of D . The third term, $\text{Bias}^2(q_D, Y)$ is the distance in terms of loss between the expected classifier and the expected label. This can be thought of as a measure of precision of our learning algorithm.

Note that we developed two things: On the one hand, we derived quantities that measure the notions of bias and variance. On the other hand, by virtue of these quantities appearing in the error decomposition 2.1.1, we have quantified the *effect* of these quantities have on the prediction error. This distinction is often overlooked because for many commonly used losses the quantities and their effects coincide. However, in general this is not necessarily true. We are particularly interested in the classification zero-one loss and a decomposition of it that is independent of the target has been proven to not exist **[todo]**. We will see that, while we cannot give a bias-variance decomposition for any loss, we can give a more general bias-variance-*effect* decomposition for which the former is a special case for some losses.

The variance-effect is the expected change in loss caused by using q_D instead of the non-random centroid q^* . Likewise, the bias-effect is the expected change in loss caused by using the expected model instead of the expected label.² Formally:

$$\begin{aligned}\text{Variance-Effect} &=_{\text{def}} \mathbb{E}_{D,Y} [L(Y, q_D) - L(Y, q^*)] \\ \text{Bias-Effect} &=_{\text{def}} \mathbb{E}_{D,Y} [L(Y, q^*) - L(Y, y^*)]\end{aligned}$$

Putting this together, we can state a generalised bias-variance decomposition. The classical bias-variance decomposition for squared error (2.1.1) is a special case of this (??).

Theorem 4.2.1 (*Bias-Variance-Effect-Decomposition*) For any loss function L , it holds that

$$\mathbb{E}_{(X,Y),D} [L(Y, q_D)] = \underbrace{\mathbb{E}_Y [L(Y, y^*)]}_{\text{noise}} + \underbrace{\mathbb{E}_{(X,Y)} [L(Y, q^*) - L(Y, y^*)]}_{\text{bias-effect}} + \underbrace{\mathbb{E}_{(D,Y)} [L(Y, q_D) - L(Y, q^*)]}_{\text{variance-effect}}$$

This decomposition holds for *any* loss function L since, by linearity of expectation, the individual terms on the right-hand-side simply cancel out.

2: In the original publication **[todo]**, bias-effect is called the *systematic effect*, i.e. the effect of the systematic components. However, it is clearer to call this *bias-effect*, particularly when we begin to introduce notions of diversity in 2.2.1.

4.3. Diversity and Diversity-Effect

chapter introduction Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

4.3.1. Measures of ensemble diversity

Maybe a table with all these measures? Would look nice. But probably not too relevant. Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

4.3.2. Ambiguity

As we have seen in 2.2.1, formally expressing the notion of ensemble diversity is not straightforward.

We take inspiration from 2.1 and derive a measure of ensemble diversity by considering the *effect* of diversity on the ensemble error. If we consider the members to be constructed according to a parameter Θ , a reasonable measure of the member performance is its loss in expectation over the parameter distribution: $\mathbb{E}_{\Theta} [L(Y, q(X; \Theta))] \approx \frac{1}{M} \sum_{i=1}^M L(Y, q(X; \Theta_i))$. What can we hope to gain from employing an ensemble \bar{q} , which combines the output of individual members q_1, \dots, q_M over just using a single member model? The *ensemble improvement* for finite ensembles is

$$\frac{1}{M} \sum_{i=1}^M L(Y, q_i) - L(Y, \bar{q})$$

Conveniently, this is also the *effect of ensembling* on the error. Due to this, we will refer to this quantity as *ambiguity-effect*.

Theorem 4.3.1 (*Ambiguity-Effect decomposition [todo]*) For any loss function L , target label Y , ensemble members q_1, \dots, q_M with combiner \bar{q}

$$L(Y, \bar{q}) = \frac{1}{M} \sum_{i=1}^M L(Y, q_i) - \underbrace{\left(\frac{1}{M} \sum_{i=1}^M L(Y, q_i) - L(Y, \bar{q}) \right)}_{\text{Ambiguity-Effect / Ensemble Improvement}}$$

At this point, it is not yet clear whether the ensemble improvement is even nonnegative, i.e. that ensembling does not hurt performance. We will address this in ??.

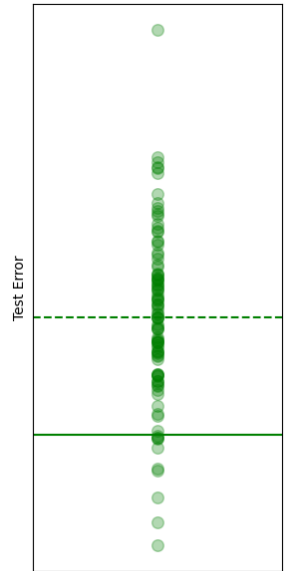


Figure 4.3.: The spread of individual tree predictions in a random forest ensemble. Glyphs correspond to test errors of individual trees. The dashed line is the average test error of individual trees $\frac{1}{M} \sum_{i=1}^M L(y, q_i)$. The solid line is the test error of the ensemble $L(y, \bar{q})$. The difference between these values is the ensemble improvement or ambiguity-effect.

Similar to variance, ambiguity and ambiguity-effect are measures of spread. Variance measures the spread of training error across models trained with different draws of the training dataset D around a model that is a centroid with respect to the distribution of D . Similarly, ambiguity measures the spread of individual member model errors around a model that is centroid with respect to the distribution of Θ , namely the combiner \bar{q} . In case of squared loss, this is indeed the statistical variance. For other losses, this is a different quantity.

4.3.3. The Bias-Variance-Diversity-Effect decomposition

Theorem 4.3.2 (*Bias-Variance-Diversity decomposition*)

4.4. Bregman Divergences

Note that bias-effect, variance-effect and ambiguity-effect are all of the form $\mathbb{E} [L(Y, \circ) - L(Y, \square)]$ and depend directly on the target label Y . With variance-effect, we have captured the *effect* of variations between different training datasets) on the prediction error. With ambiguity-effect, we have captured the effect of variations between the different member models on the prediction error. We have already seen that for the squared error, the variance-effect coincides with familiar notion of statistical variance between the predictions. This is convenient, since we can then estimate variance (and its effect) independently of the target label Y . This means that variance can be compared across different datasets. The same applies to ambiguity³. Also, labelled data can be scarce in practise, while examples alone may be easier to obtain.

3: Bias will always be indirectly dependent on the target labels via its systematic parts.

We will now define a class of losses for which the effects reduce to the quantities themselves. This class covers many widely used loss functions and thus allows us to formulate a unified bias-variance-decomposition for all of these. However, for some other losses – in particular the 0-1-loss for classification – such a decomposition is proven to not exist (see [todo]). In these cases, we can still consider the more general effect decompositions (see ??).

Definition 4.4.1 (*Bregman Divergence [todo]*) The Bregman divergence $B_\phi(\mathbf{p}, \mathbf{q})$ is defined based on a generator function ϕ as follows:

$$B_\phi(\mathbf{p}, \mathbf{q}) := \phi(\mathbf{p}) - \phi(\mathbf{q}) - \langle \nabla \phi(\mathbf{q}), (\mathbf{p} - \mathbf{q}) \rangle$$

where $\langle \cdot, \cdot \rangle$ is the inner product, $\nabla \phi(\mathbf{q})$ is the gradient vector of ϕ at \mathbf{q} and $\phi : \mathcal{S} \rightarrow \mathbb{R}$ is a strictly convex function on a convex set $\mathcal{S} \subseteq \mathbb{R}^k$ such that it is differentiable on the relative interior of \mathcal{S} .

Bregman divergences have the one key property that their effect terms collapse. This is given by the following two lemmas.

Lemma 4.4.1 ([pfau], Theorem 0.1 (b)) Let the generator $\phi : \mathcal{S} \rightarrow \mathbb{R}$ be a strictly convex, differentiable function. Let Y be a random variable on \mathcal{S} and $\mathbf{x}^\star = \arg \min_{\mathbf{z}} \mathbb{E} [B_\phi(\mathbf{Y}, \mathbf{z})] = \mathbb{E} [Y]$. Then, for any $\mathbf{q} \in \mathcal{S}$, it holds that

$$B_\phi(\mathbf{y}^\star, \mathbf{q}) = \mathbb{E} [B_\phi(\mathbf{Y}, \mathbf{q})] - \mathbb{E} [B_\phi(\mathbf{Y}, \mathbf{y}^\star)]$$

where the expectation is over the distribution of Y .

This shows that bias-effect collapses to bias for Bregman divergences: $B_\phi(\mathbf{y}^\star, \mathbf{q}^\star) = \mathbb{E} [B_\phi(\mathbf{Y}, \mathbf{q}^\star) - B_\phi(\mathbf{Y}, \mathbf{y}^\star)]$

Lemma 4.4.2 (Generalised from [ref:wood23], [pfau_GeneralizedBiasVarianceDecomposition_])

Let \mathbf{y} be a random vector. Let \mathbf{q}_Z be a random vector dependent on another random variable Z . Then it holds that

$$B_\phi(\mathbf{q}^\star, \mathbf{q}) = \mathbb{E}_Y [B_\phi(\mathbf{Y}, \mathbf{q}) - B_\phi(\mathbf{Y}, \mathbf{q}^\star)]$$

$$\text{if } \mathbf{q}^\star = (\nabla\phi)^{-1}(\mathbb{E}_Z [\nabla\phi(\mathbf{q})])$$

Proof.

$$\mathbb{E}_{Y,Z} [B_\phi(\mathbf{y}, \mathbf{q}) - B_\phi(\mathbf{y}, \mathbf{q}^\star)] = \dots$$

□

This shows that variance- and ambiguity-effect collapse to variance and ambiguity:

$$B_\phi(\mathbf{q}^\star, \mathbf{q}) = \mathbb{E} [B_\phi(\mathbf{Y}, \mathbf{q}^\star) - B_\phi(\mathbf{Y}, \mathbf{q})]$$

$$B_\phi(\mathbf{q}^\star, \mathbf{q}_\Theta) = \mathbb{E} [B_\phi(\mathbf{Y}, \mathbf{q}^\star) - B_\phi(\mathbf{Y}, \mathbf{q})] \quad \text{for } \mathbf{q}^\star = \mathbb{E}_\Theta [\mathbf{q}_\Theta]$$

This yields a generalised bias-variance decomposition for bregman divergences as a special case of the bias-variance-effect decomposition (2.1.2).

Theorem 4.4.3 (Bias-variance-decomposition for Bregman divergences)

$$\mathbb{E}_{(X,Y),D} [B_\phi(\mathbf{Y}, \mathbf{q})] = \underbrace{\mathbb{E}_{Y|X} [B_\phi(\mathbf{Y}, \mathbf{y}^\star)]}_{\text{noise}} + \underbrace{B_\phi(\mathbf{y}^\star, \mathbf{q}^\star)}_{\text{bias}} + \underbrace{\mathbb{E}_D [B_\phi(\mathbf{q}^\star, \mathbf{q})]}_{\text{variance}}$$

Theorem 4.4.4 (Bias-Variance-Diversity decomposition)

4.5. Bias, Variance and Covariance

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

5. Growth Strategies informed by Diversity

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

5.1. Dynamic Random Forests

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

5.2. Generalized Negative Correlation Learning

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like "Huardest gefburn"? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

6. Conclusion

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.

A. Appendix

APPENDIX

B. Some more blindtext

Hello, here is some text without a meaning. This text should show what a printed text will look like at this place. If you read this text, you will get no information. Really? Is there no information? Is there a difference between this text and some nonsense like “Huardest gefburn”? Kjift – not at all! A blind text like this gives you information about the selected font, how the letters are written and an impression of the look. This text should contain all letters of the alphabet and it should be written in of the original language. There is no need for special content, but the length of words should match the language.