# Abstract

Each word or token is treated as one data, obtain the dense vector of word or token are representing as a sparse vector (the dimensional is 2,035,523) and use classification algorithms to do it such as Gaussian Process.
Because the dimension is too high, PCA is used to reduce dimension (the dimensional is 3000).
The value of ER is: 0.43
The value of MNLP is: 42

# Introduction

This is a problem of sequence annotation, and it is also a classification problem. Here we consider it as a classification problem, each word corresponds to a category, regardless of the context. Of course, this accuracy rate will be very low. We use Gaussian Process to classification.
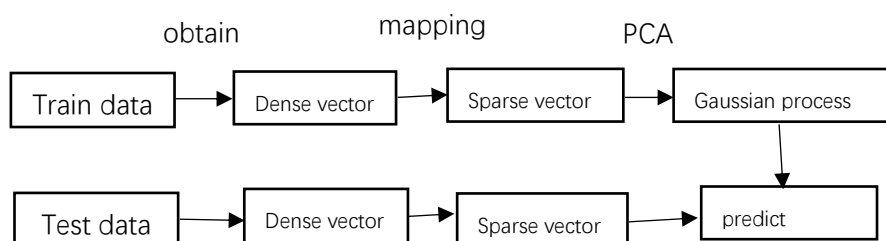
# Model

$$\text{Data:} D = \{(X_i, Y_i)\}_i^N, X \in R^D, y \in \{0,1,\cdots 22\} \quad \text{-> there are 23 categories.}$$

$$\text{Input:} (X)_{D*N}, \quad Target: (y)_{N*1}$$

$$\text{Prior:} f \sim GP\big(0, K(x, \acute{x}; \theta)\big), p(f|X) = N(f|0, K)$$

$$\text{Likelihood:} p(y|f) = \sigma\big(f(x)\big) = \prod_{i=1}^{N} p(y_i|f(x_i))$$

$$\sigma(z) = \frac{1}{1 + \exp(z)}$$

# Inference

Compute predictive distribution of latent functions:

$$p(f_*|X, y, x_*) = \int p(f_*|X, x_*, f) \cdot p(f|X, y) df$$

$$p(f|X, y) \propto p(f|X) \cdot p(y|f)$$

Compute probabilistic predictions:

$$p(y_* = +1|X, y, x_*) = \int \sigma(f_*) \cdot p(f_*|X, y, x_*) df_*$$

# Parameter Estimation

Use Laplace approximation

$$p(f|X, y, \theta) \approx N(f|\hat{f}, A^{-1})$$

where: $\hat{f} = \text{argmax}\, p(y|f, \theta) \cdot p(f|X, \theta)$ and A is the Hessian of the negative log-posterior

evaluated at $\hat{f}$.

# Results

There are two indicators for evaluating the performance of models. ER and MNLP

$$\text{ER} = 1 - \frac{1}{N} \sum_{i=1}^{N} \frac{1}{T_i} \sum_{k=1}^{T_i} [\![\hat{y}_k^{(i)} = y_i^{(i)}]\!]$$

$$\text{MNLP} = -\frac{1}{N} \sum_{i=1}^{N} \frac{1}{T_i} \sum_{k=1}^{T_i} \sum_{j=1}^{C} [\![y_k^{(i)} = y]\!] log P_{model}(y_k^{(i)}|x^{(i)})$$

The value of ER is: 0.43
The value of MNLP is: 42
If we use the method of cross-validation, the performance will be better.

# References

Gaussian Process Classification, Approximations and Other GP Models - UNSW Sydney
http://scikit-learn.org
Variational Learning of GP Models - UNSW Sydney

# Appendix

The dimension of the feature space is not equal 2035523, only have 90000 dimension, so we using the value of 90000 to mapping sparse vector, will be faster and use less memory.