# Quick Permutation Test (QuiPT)

Piotr Sobczyk[1], Michał Burdukiewicz[2]

[1]Wrocław University of Technology, Institute of Mathematics and Computer Science, Poland

[2]University of Wrocław, Department of Genomics, Poland

## Outline

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

n-gram definition
Positioned n-grams

# Outline

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

n-gram definition
Positioned n-grams

n-grams (k-tuples) are sets of n characters derived from the input sequence(s). They may form continuous sub-sequences or be discontinuous.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

n-gram definition
Positioned n-grams

|   | X1 | X2 | X3 | X4 | X5 | X6 |
|---|----|----|----|----|----|----|
| 1 | 4  | 1  | 4  | 1  | 3  | 3  |
| 2 | 1  | 4  | 3  | 1  | 2  | 2  |
| 3 | 1  | 4  | 1  | 1  | 1  | 1  |

Sample sequences.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

n-gram definition
Positioned n-grams

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 2 | 0 | 2 | 2 |
| 2 | 2 | 1 | 1 |
| 5 | 0 | 0 | 1 |

Unigrams.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

n-gram definition
Positioned n-grams

| X1_1_0 | X2_1_0 | X3_1_0 | X4_1_0 | X5_1_0 | X6_1_0 | X1_2_0 |
|:------:|:------:|:------:|:------:|:------:|:------:|:------:|
| 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 | 1 | 1 | 0 |

Part of possible unigrams with position information.

Positioned n-gram data is binary.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

n-gram definition
Positioned n-grams

Number of possible positioned n-grams:

$$n_{max} = L \times m^n$$

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

n-gram definition
Positioned n-grams

n-grams
**Permutation test**
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

# Outline

n-grams
**Permutation test**
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

1. Calculate test statistic for the given positioned n-gram and etiquettes ($T_R$).

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

1. Calculate test statistic for the given positioned n-gram and etiquettes ($T_R$).
2. Permute counts of n-grams and calculate permuted test statistic ($T_P$).

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

1. Calculate test statistic for the given positioned n-gram and etiquettes ($T_R$).
2. Permute counts of n-grams and calculate permuted test statistic ($T_P$).
3. Repeat step 2. N times.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

1. Calculate test statistic for the given positioned n-gram and etiquettes ($T_R$).
2. Permute counts of n-grams and calculate permuted test statistic ($T_P$).
3. Repeat step 2. N times.
4. Calculate p-value using:

$$\text{p-value} = \frac{N_{T_P > T_R}}{N}$$

$N_{T_P > T_R}$ is number of times when $T_P$ was bigger than $T_R$

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

- Model independent.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

- Model independent.
- Statistic independent.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

- Computationally expensive (number of cases, number of features).

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

- Computationally expensive (number of cases, number of features).
- Single feature analysis (no feature interaction).

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Advantages
Drawbacks
p-value resolution

The number of permutations is inversely proportional to the resolution of p-value.

Example: with $10\times10^6$ permutation the smallest possible p-values are: 0, $1\times10^{-6}$, $2\times10^{-6}$ and so on.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

# Outline

1. n-grams
   - n-gram definition
   - Positioned n-grams

2. Permutation test
   - Advantages
   - Drawbacks
   - p-value resolution

3. QuiPT
   - Contingency tables
   - Multinomial distribution of target-feature relationship
   - Advantages over permutation test

4. Simulation scheme
   - Power of the test
   - False significant features

5. Conclusion
   - Advantages over permutation test

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

The binary positioned n-gram data tabulated by binary label can be easily described in 2d contingency table.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

| sequence ID | feature | target |
|:-----------:|:-------:|:------:|
| 1 | 1 | 0 |
| 2 | 1 | 0 |
| 3 | 0 | 0 |
| 4 | 1 | 1 |
| 5 | 0 | 1 |
| . . . | . . . | . . . |

Positioned n-grams with a label.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

|   | target | feature |
|---|--------|---------|
| 0 | $n_{1,1}$ | $n_{1,0}$ |
| 1 | $n_{0,1}$ | $n_{0,0}$ |

Contingency table.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

Test statistics used by QuiPT (information gain, Kullback-Leibler divergence) measure inbalance of contingency tables.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

If probability that target equals 1 is $p$ and probability that feature equals 1 is $q$ and feature and target are independent then each of them has the following probabilities

$$P(Target, Feature) = (1, 1)) = p \cdot q$$

$$P(Target, Feature) = (1, 0)) = p \cdot (1 - q)$$

$$P(Target, Feature) = (0, 1)) = (1 - p) \cdot q$$

$$P(Target, Feature) = (0, 0)) = (1 - p) \cdot (1 - q)$$

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

$$F(n_{1,1}, n_{1,0}, n_{0,1}, n_{0,0}) = \binom{n}{n_{1,1}} (p \cdot q)^{n_{1,1}} \binom{n - n_{1,1}}{n_{1,0}} (p \cdot (1 - q))^{n_{1,0}}$$
$$\binom{n - n_{1,1} - n_{1,0}}{n_{0,1}} ((1 - p) \cdot q)^{n_{0,1}}$$
$$\binom{n - n_{1,1} - n_{1,0} - n_{0,1}}{n_{0,0}} ((1 - p) \cdot (1 - q))^{n_{0,0}}$$

In addition to this: $n_{1,\cdot} = n_{1,1} + n_{1,0}$ and $n_{\cdot,1} = n_{1,1} + n_{0,1}$ are known and fixed.
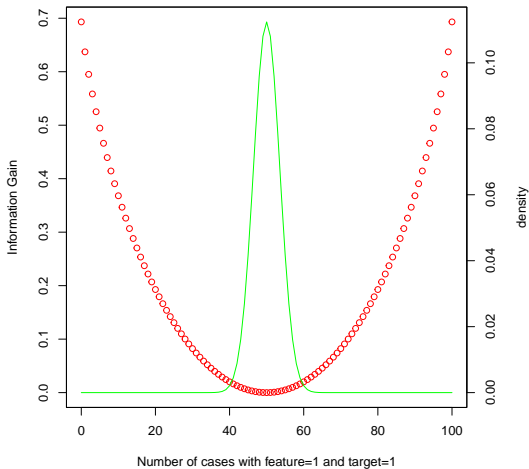
n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

- $n_{1,1}$ is from range $[0, min(n_{\cdot,1}, n_{1,\cdot})]$.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

- $n_{1,1}$ is from range $[0, min(n_{\cdot,1}, n_{1,\cdot})]$.
- The probability of certain contingency table is given as the conditional distribution, as impose restrictions on two parameters $n_{\cdot,1}$ and $n_{1,\cdot}$.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

- $n_{1,1}$ is from range $[0, min(n_{\cdot,1}, n_{1,\cdot})]$.
- The probability of certain contingency table is given as the conditional distribution, as impose restrictions on two parameters $n_{\cdot,1}$ and $n_{1,\cdot}$.
- The test statistic is computed for each possible value of $n_{1,1}$.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

- $n_{1,1}$ is from range $[0, min(n_{\cdot,1}, n_{1,\cdot})]$.
- The probability of certain contingency table is given as the conditional distribution, as impose restrictions on two parameters $n_{\cdot,1}$ and $n_{1,\cdot}$.
- The test statistic is computed for each possible value of $n_{1,1}$.
- The distribution of test statistics under hypothesis that target and feature are independant is computed using values from 3.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

| | Target | Feature | Freq |
|---|---|---|---|
| 1 | 0 | 0 | 50 |
| 2 | 1 | 0 | 50 |
| 3 | 0 | 1 | 50 |
| 4 | 1 | 1 | 50 |

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
Advantages over permutation test

- QuiPT is faster.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Contingency tables
Multinomial distribution of target-feature relationship
**Advantages over permutation test**

- QuiPT is faster.
- Using the exact distribution of possible values of the criterion QuiPT yields precise small p-values without increasing the computation time.

n-grams
Permutation test
QuiPT
**Simulation scheme**
Conclusion

Power of the test
False significant features

# Outline

1. n-grams
   - n-gram definition
   - Positioned n-grams

2. Permutation test
   - Advantages
   - Drawbacks
   - p-value resolution

3. QuiPT
   - Contingency tables
   - Multinomial distribution of target-feature relationship
   - Advantages over permutation test

4. **Simulation scheme**
   - Power of the test
   - False significant features

5. Conclusion
   - Advantages over permutation test

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Power of the test
False significant features

1. Random 4000 sequences (20 nucleotides each). The half of the sequences has label 0.

n-grams
Permutation test
QuiPT
Simulation scheme
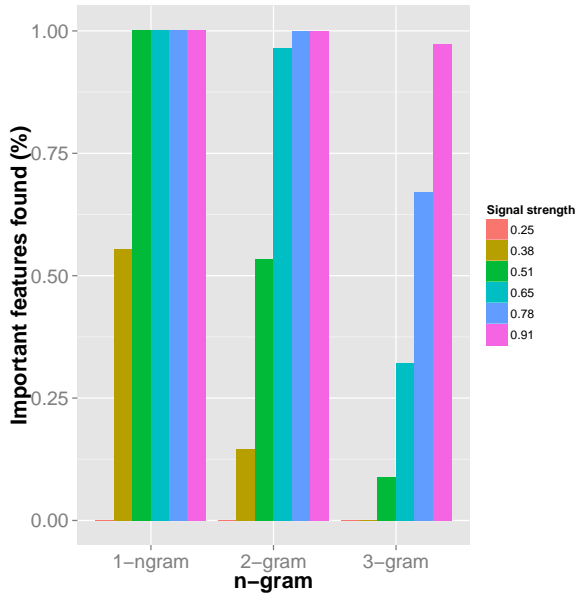Conclusion

Power of the test
False significant features

1. Random 4000 sequences (20 nucleotides each). The half of the sequences has label 0.
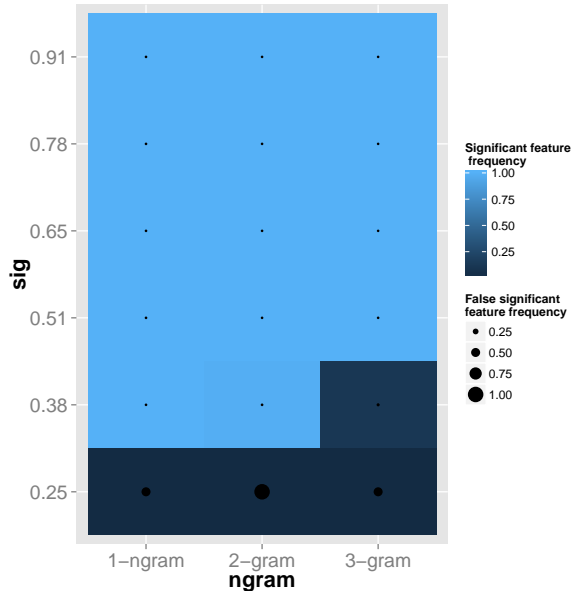2. Choose a single position between 3 and 18 (to avoid border cases).

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Power of the test
False significant features

1. Random 4000 sequences (20 nucleotides each). The half of the sequences has label 0.
2. Choose a single position between 3 and 18 (to avoid border cases).
3. Resample nucleotides at chosen position. The dominant nucleotoide has probabilitiy of occurence $p_d = 0.25$. Other nucleotides have probability of occurence $p_o = (1 - p_d)/3$.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Power of the test
False significant features

1. Random 4000 sequences (20 nucleotides each). The half of the sequences has label 0.
2. Choose a single position between 3 and 18 (to avoid border cases).
3. Resample nucleotides at chosen position. The dominant nucleotide has probabilitiy of occurence $p_d = 0.25$. Other nucleotides have probability of occurence $p_o = (1 - p_d)/3$.
4. Perform QuiPT (Information Gain) and choose significant features (with p-value $< 0.001$).

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Power of the test
False significant features

1. Random 4000 sequences (20 nucleotides each). The half of the sequences has label 0.

2. Choose a single position between 3 and 18 (to avoid border cases).

3. Resample nucleotides at chosen position. The dominant nucleotide has probabilitiy of occurence $p_d = 0.25$. Other nucleotides have probability of occurence $p_o = (1 - p_d)/3$.

4. Perform QuiPT (Information Gain) and choose significant features (with p-value $< 0.001$).

5. Iterate steps 1-4 over other values of $p_d$ - 0.38, 0.51, 0.65, 0.78, 0.91.

n-grams
Permutation test
QuiPT
Simulation scheme
Conclusion

Power of the test
False significant features

1. Random 4000 sequences (20 nucleotides each). The half of the sequences has label 0.
2. Choose a single position between 3 and 18 (to avoid border cases).
3. Resample nucleotides at chosen position. The dominant nucleotide has probabilitiy of occurence $p_d = 0.25$. Other nucleotides have probability of occurence $p_o = (1 - p_d)/3$.
4. Perform QuiPT (Information Gain) and choose significant features (with p-value $< 0.001$).
5. Iterate steps 1-4 over other values of $p_d$ - 0.38, 0.51, 0.65, 0.78, 0.91.
6. Repeat steps 1-5 200 times.

n-grams
Permutation test
QuiPT
**Simulation scheme**
Conclusion

Power of the test
False significant features

n-grams
Permutation test
QuiPT
Simulation scheme
**Conclusion**

Advantages over permutation test

# Outline

1. n-grams
   - n-gram definition
   - Positioned n-grams

2. Permutation test
   - Advantages
   - Drawbacks
   - p-value resolution

3. QuiPT
   - Contingency tables
   - Multinomial distribution of target-feature relationship
   - Advantages over permutation test

4. Simulation scheme
   - Power of the test
   - False significant features

5. Conclusion
   - Advantages over permutation test

n-grams
Permutation test
QuiPT
Simulation scheme
**Conclusion**

Advantages over permutation test

Quick permutation test is a powerful and quick equivalent of permutation test in binary feature-binary target testing scenario.